

String Inference from Longest-Common-Prefix Array

Juha Kärkkäinen¹, Marcin Piatkowski², and Simon J. Puglisi¹

¹ Helsinki Institute of Information Technology (HIIT) and
Department of Computer Science, University of Helsinki, Finland

² Faculty of Mathematics and Computer Science,
Nicolaus Copernicus University, Toruń, Poland

{juha.karkkainen,simon.puglisi}@cs.helsinki.fi
marcin.piatkowski@mat.umk.pl

Abstract. The suffix array, perhaps the most important data structure in modern string processing, is often augmented with the longest common prefix (LCP) array which stores the lengths of the LCPs for lexicographically adjacent suffixes of a string. Together the two arrays are roughly equivalent to the suffix tree with the LCP array representing the tree shape.

In order to better understand the combinatorics of LCP arrays, we consider the problem of inferring a string from an LCP array, i.e., determining whether a given array of integers is a valid LCP array, and if it is, reconstructing some string or all strings with that LCP array. There are recent studies of inferring a string from a suffix tree shape but using significantly more information (in the form of suffix links) than is available in the LCP array.

We provide two main results. (1) We describe two algorithms for inferring strings from an LCP array when we allow a generalized form of LCP array defined for a multiset of cyclic strings: a linear time algorithm for binary alphabet and a general algorithm with polynomial time complexity for a constant alphabet size. (2) We prove that determining whether a given integer array is a valid LCP array is NP-complete when we require more restricted forms of LCP array defined for a single cyclic or non-cyclic string or a multiset of non-cyclic strings. The result holds whether or not the alphabet is restricted to be binary. In combination, the two results show that the generalized form of LCP array for a multiset of cyclic strings is fundamentally different from the other more restricted forms.

Keywords: LCP array, string inference, BWT, suffix array, suffix tree, NP-hardness

1 Introduction

For a string X of n symbols, the suffix array (SA) [22] contains pointers to the suffixes of X , sorted in lexicographical order. The suffix array is often augmented with a second array — the longest common prefix (LCP) array — storing the length of the longest common prefix between lexicographically adjacent suffixes; i.e., $\text{LCP}[i]$ is the length of the LCP of suffixes $X[\text{SA}[i]..n]$ and $X[\text{SA}[i-1]..n]$. The two arrays are closely connected to the suffix tree [31] — the compacted trie of all the string’s suffixes: the entries of SA correspond to the leaves of the suffix tree, and the LCP array entries tell the string depths of the lowest common ancestors of adjacent leaves, defining the shape of the tree (see Fig. 2 in the appendix). For decades these data structures have been central to string processing; see [4] for a history and an overview, and [1,3,15,29,25] for further details on myriad applications.

Given both the suffix and the LCP array, the corresponding string is unique up to renaming of the characters and is easy to reconstruct: zeros in the LCP array tell where the first character changes in the lexicographical list of the suffixes, and the suffix array tells how to permute those first characters to obtain the string. Given just the suffix array, we can easily reconstruct a corresponding string where all characters are different, and it is not difficult to characterize strings with a given suffix array [5,27,21]. In essence, the suffix array determines a set of positions in the LCP array that must be zero. Specifically, for any i let j and k be integers such that $\text{SA}[j] = \text{SA}[i-1] + 1$ and $\text{SA}[k] = \text{SA}[i] + 1$. Then, if $k < j$, we must have $\text{LCP}[i] = 0$. For any

other position, we can freely and independently decide whether the value is zero or not, and as described above, the zero positions together with the suffix array determine the string.

In this paper, we consider the problem of similarly reconstructing strings from an LCP array without the suffix array. As mentioned above, the LCP array determines the shape of the suffix tree, i.e., the suffix tree without edge or leaf labels. Notice that the LCP array specifies the label lengths for internal edges but not for leaf edges, which would allow trivial inference of the suffix array. String inference from the suffix tree shape has recently been considered by three different sets of authors [19,6,30]. However, all of them assume that the suffix tree is augmented with significant additional information, namely *suffix links* (see Fig. 2), which makes the task much easier. Indeed, our new algorithms essentially reconstruct suffix links from the LCP array. According to Cazaux and Rivals [6], the case without suffix links was considered but not solved in [26]. We are also aware that others have considered it but without success [2].

To fully define the problem, we have to specify what kind of strings we are trying to infer. Often suffix trees and suffix arrays are defined for *terminated strings* that are assumed to end with a special symbol \$ that is different from and lexicographically smaller than any other symbol. The alternative is an *open-ended string* where no assumption is made on the last symbol. For suffix and LCP arrays the only change from omitting the terminator symbol is dropping the first element (which is always zero in the LCP array), but the suffix tree can change considerably because some suffixes can be prefixes of other suffixes and thus are not represented by a leaf (see Fig. 3). Inferring open-ended strings from a suffix tree (with suffix links) is studied by Starikovskaya and Vildhøj [30], who show that any string can be appended by additional characters without changing the suffix tree shape (thus the term open-ended). However, such an extension can change the suffix and LCP arrays a great deal (see Fig. 4), i.e., with the arrays a string is never truly open-ended but has at least an implicit terminator.

To get rid of even an implicit terminator, we consider a third type of strings, *cyclic strings*, where we use rotations in place of suffixes (see Figs. 5–7). For a terminated string, replacing suffixes with rotations causes no changes to the suffix/rotation array or the LCP array. Thus any integer array that is a valid LCP array for a terminated string is always a valid LCP array for a cyclic string too, but the opposite is not true. For example, the LCP array for the cyclic string *aababa* is $(2, 1, 3, 0, 2)$, which is not a valid LCP array for any non-cyclic string. In this sense, the cyclic string case is strictly more general. An even more striking example is a non-primitive string, such as *abab*, that has two or more identical rotations. For reasons explained below, instead of rotations we use *cyclic suffixes* which are infinite repetitions of rotations. Thus the LCP array for the cyclic string *abab* is $(\omega, 0, \omega)$, where ω denotes the positions of two adjacent identical cyclic suffixes.

Finally, we may have a joint suffix array for a collection of strings, where we have all suffixes of all strings in lexicographical order, and the corresponding LCP array. In the terminated version, each string is terminated with a distinct terminator symbol. If we have an LCP array for a collection of open-ended strings, adding the terminator symbols simply prepends one zero for each terminator. The LCP array for a collection of terminated strings is identical to the LCP array of the concatenation of the strings. Thus the generalization from single strings to string sets does not add to the set of valid LCP arrays for terminated strings, but it does for cyclic strings. For example the LCP array for a string set $\{aa, b\}$ is $(\omega, 0)$, which is not a valid LCP array for any single string. For multiple cyclic strings, it is important to use cyclic suffixes instead of rotations because the result can be different (e.g., the set $\{ab, aba\}$).

Now we are ready to formally define the problem of String Inference from LCP Array (SILA). In the decision version, we are given an array of integers (and possibly ω 's) and asked if the array is a valid LCP array of some string. If the answer is yes, the reporting version may also output some such string, and possibly a characterization of all such strings. Different variants

are identified by a prefix: S for a string set; T, O, or C for terminated, open-ended or cyclic; and B for a binary alphabet (where terminators are not counted). For example, BCSSILA stands for Binary Cyclic String Set Inference from LCP Array. As discussed above, and summarized in the following result (with a proof in the appendix), the non-cyclic variants are essentially equivalent, but the cyclic variants are more general.

Proposition 1. *There are polynomial time reductions from BTSILA to BOSILA, BTSSILA, BOSSILA, TSILA, OSILA, TSSILA, and OSSILA.*

Our Contribution. Our first result is a linear time algorithm for BCSSILA. For a valid LCP array the algorithm outputs a string, which is the Burrows-Wheeler transform (BWT) of the solution string set. This relies on a generalization of the BWT for multisets of cyclic strings developed in [23,20]. There can be more than one multiset of strings with the same BWT but the class of such string collections is simple and well characterized in [20]. The algorithm also outputs a set of substring swaps such that applying any combination of the swaps on the BWT produces another BWT of a solution, and any BWT of a solution can be produced by such a combination of swaps. Thus we have a complete characterization of all solutions. The number of swaps can be linear and thus the number of distinct solutions can be exponential. We also present an algorithm for CSSILA, i.e., without a restriction on the alphabet size, that has a polynomial time complexity for any constant alphabet size.

Our second result is a proof, by a reduction from 3SAT, that (the decision version of) BCSILA, and thus CSILA, is NP complete. Therefore, even though the BCSSILA algorithm produces a characterization of all solutions, it is NP hard to determine whether one of the solutions is a single string. Furthermore, we modify the reduction to prove that BTSILA is NP complete too. By Proposition 1, this shows that all variants of SILA mentioned above except (B)CSSILA are NP complete. Since CSSILA is in P for constant alphabet sizes, this leaves the complexity of CSSILA for larger alphabets as an open problem.

Related Work. String inference from partial information is a classic problem in string processing, dating back some 40 years to the work of Simon [28], where reconstructing a string from a set of its subsequences is considered. Since then, string inference from a variety of data structures has received a considerable amount attention, with authors considering border arrays [12,11,10], parameterized border arrays [18], the Lyndon factorization [24], suffix arrays [5,21], KMP failure tables [11,13], prefix tables [7], cover arrays [9], and directed acyclic word graphs [5]. The motivation for studying most string inference problems is to gain a deeper understanding of the combinatorics of the data structures involved, in order to design more efficient algorithms for their construction and use.

A (somewhat tangentially) related result to ours is due to He et al. [16], who prove that it is NP hard to infer a string from the longest-previous-factor (LPF) array. It is well known that LPF is a permutation of LCP [8] but otherwise it is a quite different data structure. For example, it is in no way concerned with lexicographical ordering. Like our NP-hardness proof, He et al.’s reduction is from 3-SAT, but the details of each reduction appear to be very different. Moreover, their construction requires an unbounded alphabet while our construction works for a binary alphabet and thus for any alphabet.

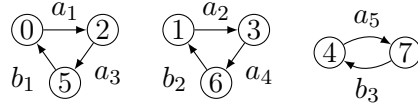
To the best of our knowledge, all of the previous string inference problems aim at obtaining a single non-cyclic string from some data structure, and we are the first to consider the generalizations to cyclic strings and to string sets, and as our results show, this makes a crucial difference. As explained in the next section, the generalizations arise naturally from the generalized BWT introduced in [23], which also played a central role in another recent result on the combinatorics of LCP arrays [20].

2 Basic notions

Let v be a string of length n and let \hat{v} be obtained from v by sorting its characters. The *standard permutation* [14,17] of v is the mapping $\Psi_v : [0..n) \rightarrow [0..n)$ such that for every $i \in [0..n)$ it holds $\hat{v}[i] = v[\Psi_v(i)]$ and for any $\hat{v}[i] = \hat{v}[j]$ the relation $i < j$ implies $\Psi_v(i) < \Psi_v(j)$. In other words, Ψ_v corresponds to the stable sorting of the characters. Let $C = \{c_i\}_{i=1}^s$ be the disjoint cycle decomposition of Ψ_v . We define the inverse Burrows–Wheeler transform IBWT as the mapping from v into a multiset of cyclic strings $W = \{\{w_i\}_{i=1}^s\}$ such that for any $i \in [1..s]$ and $j \in [0..|c_i|)$, $w_i[j] = v[\Psi_v(c_i[j])]$.

Example 1. For $v = bbaabaaa$, we have $\text{IBWT}(v) = \{\{aab, aab, ab\}\}$ as illustrated in the following table (showing \hat{v} and Ψ_v) and figure (showing the cycles of Ψ_v as a graph). The character subscripts are provided to make it easier to ensure stability.

| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|--------------|-------|-------|-------|-------|-------|-------|-------|-------|
| $v[i]$ | b_1 | b_2 | a_1 | a_2 | b_3 | a_3 | a_4 | a_5 |
| $\hat{v}[i]$ | a_1 | a_2 | a_3 | a_4 | a_5 | b_1 | b_2 | b_3 |
| $\Psi_v[i]$ | 2 | 3 | 5 | 6 | 7 | 0 | 1 | 4 |



The elements of W are primitive cyclic strings. *Cyclic* means that all rotations of a string are considered equal. For example, aab , aba and baa are all equal. A string is *primitive* if it is not a concatenation of multiple copies of the same string. For example, aab is primitive but $aabaab$ is not. For any alphabet Σ , the mapping IBWT is a bijection between the set Σ^* of all (non-cyclic) strings and the multisets of primitive cyclic strings over Σ [23].

The set of positions of W is defined as the set of integer pairs $\text{pos}(W) := \{\langle i, p \rangle : i \in [1..s], p \in [0..|w_i|)\}$. For a position $\langle i, p \rangle \in \text{pos}(W)$ we define a *cyclic suffix* $W_{\langle i, p \rangle}$ as the infinite string that starts at $\langle i, p \rangle$, i.e., $W_{\langle i, p \rangle} = w_i[p]w_i[p+1 \bmod |w_i|]w_i[p+2 \bmod |w_i|], \dots$. The multiset of all cyclic suffixes of W is defined as $\text{suf}(W) := \{\{W_{\langle i, p \rangle} : \langle i, p \rangle \in \text{pos}(W)\}\}$. We say that a string x occurs at position $\langle i, p \rangle$ in W if x is a prefix of the suffix $W_{\langle i, p \rangle}$.

The *(cyclic) suffix array* of a multiset of strings W is defined as an array $\text{SA}_W[j] = \langle i_j, p_j \rangle$, where $\langle i_j, p_j \rangle \in \text{pos}(W)$ for all $j \in [0..n)$ and $W_{\langle i_{j-1}, p_{j-1} \rangle} \leq W_{\langle i_j, p_j \rangle}$ for all $j \in [1..n)$. The *Burrows–Wheeler transform* (BWT) is a mapping from W into the string v defined as $v[j] = w_i[p-1 \bmod |w_i|]$, where $\langle i, p \rangle = \text{SA}_W[j]$, i.e., $v[j]$ is the character preceding the beginning of the suffix $W_{\text{SA}_W[j]}$. The BWT is the inverse of IBWT [23,20].

The *longest-common-prefix array* $\text{LCP}_W[1..n)$ is defined as $\text{LCP}_W[j] = \text{lcp}(W_{\text{SA}_W[j-1]}, W_{\text{SA}_W[j]})$ for $0 < j < n$, where $\text{lcp}(x, y)$ is the length of the longest common prefix between the strings x and y .

Example 2. For $W = \{\{ab, aab, aab\}\}$ we have

$$\begin{aligned} \text{suf}(W) &= \{(aab)^\omega, (aab)^\omega, (aba)^\omega, (aba)^\omega, (ab)^\omega, (baa)^\omega, (baa)^\omega, (ba)^\omega\} \\ \text{SA}_W &= [\langle 2, 0 \rangle, \langle 3, 0 \rangle, \langle 2, 1 \rangle, \langle 3, 1 \rangle, \langle 1, 0 \rangle, \langle 2, 2 \rangle, \langle 3, 2 \rangle, \langle 1, 1 \rangle] \\ \text{LCP}_W &= [\omega, 1, \omega, 3, 0, \omega, 2]. \end{aligned}$$

The suffixes represented by the suffix array entries can also be expressed as follows.

Lemma 1. For $i \in [0..n)$, $W_{\text{SA}_W[i]} = \hat{v}[i] \cdot \hat{v}[\Psi_v(i)] \cdot \hat{v}[\Psi_v^2(i)] \cdot \hat{v}[\Psi_v^3(i)] \dots$

2.1 Intervals.

Many algorithms on suffix arrays and LCP arrays are based on iterating over a specific types of array intervals. Next, we define these intervals and establish their key properties. For proofs and further details, we refer to [1,25].

Let $v \in \{a, b\}^n$ and $W = IBWT(v)$. Let $SA = SA_W$ be the suffix array and $LCP = LCP_W$ the LCP array of W . Note that from now on, we will assume a binary alphabet.

Definition 1 (x -interval). *An interval $[i..j]$, $0 \leq i \leq j \leq n$, is called the x -interval ($x \in \Sigma^*$) if and only if (1) x is not a prefix of $W_{SA[i-1]}$ (or $i = 0$), (2) x is a prefix of $W_{SA[k]}$ for all $k \in [i..j]$, and (3) x is not a prefix of $W_{SA[j]}$ (or $j = n$).*

In other words, in the suffix array the x -interval $SA[i..j]$ consists of all suffixes of W with x as a prefix. Thus the size $j - i$ of the interval is the number of occurrences of x in W , which we will denote by n_x .

Definition 2 (ℓ -interval). *An interval $[i..j]$, $0 \leq i < j \leq n$, is called an ℓ -interval ($\ell \in \mathbb{N} \cup \{\omega\}$) if and only if (1) $LCP[i] < \ell$ (or $i = 0$), (2) $\min LCP[i+1..j] = \ell$ (where $\min LCP[j..j] = \omega$), and (3) $LCP[j] < \ell$ (or $j = n$).*

Lemma 2. *Every nonempty x -interval is an ℓ -interval for some (unique) $\ell \geq |x|$. Every ℓ -interval is an x -interval for some string x of length ℓ .*

Corollary 1. *If an x -interval $[i..j]$ is an ℓ -interval for $\ell > |x|$, there exists a (unique) string y of length $\ell - |x|$ such that $[i..j]$ is the xy -interval.*

Thus the ℓ -intervals represent the set of all distinct x -intervals. This and the fact that the total number of ℓ -intervals is $\mathcal{O}(n)$ are the basis of many efficient algorithms for suffix arrays, see e.g., [1,25].

3 Algorithm for BCSSILA

We are now ready to describe the algorithm for string inference from an LCP array. Given an LCP array $LCP[1..n]$, our goal is to construct a string $v \in \{a, b\}^n$ such that $LCP = LCP_{IBWT(v)}$. At first, we assume that such a string v exists, and consider later what happens if the input is not a valid LCP array.

Let $RMQ_{LCP}[i..j]$ denote the *range minimum query* over the LCP array that returns the position of the minimum element in $LCP[i..j]$, i.e., $RMQ_{LCP}[i..j] = \arg \min_{k \in [i..j]} LCP[k]$. The LCP array is preprocessed in linear time so that any RMQ can be answered in constant time (see for instance [25]). Then any x -interval can be split into two subintervals as shown in the following result.

Lemma 3. *Let $[i..j]$ be an x -interval and an ℓ -interval for $\ell < \omega$, and let $k = RMQ_{LCP}[i+1..j]$. Then, for some string y of length $\ell - |x|$, $[i..k]$ is the xya -interval and $[k..j]$ is the xyb -interval.*

This approach makes it easy to recursively enumerate all ℓ -intervals. We will also keep track of ax - and bx -intervals together with any x -interval, even if we do not know x precisely. From the intervals we can determine the numbers of occurrences, n_{ax} and n_{bx} , which are useful in the inference of v :

Lemma 4. *Let $[i..j]$ be the x -interval. Then $v[i..j]$ contains exactly n_{ax} a 's and n_{bx} b 's.*

In particular, when either n_{ax} or n_{bx} drops to zero, we have fully determined $v[i..j)$ for the x -interval $[i..j)$. In such a case, the LCP array intervals have to satisfy the following property.

Lemma 5. *Let $[i_y..j_y)$ be the y -interval for $y \in \{x, ax, bx\}$. If $n_{ax} = j_{ax} - i_{ax} = 0$, then $\text{LCP}[i_{bx} + 1..j_{bx}) = 1 + \text{LCP}[i_x + 1..j_x)$, where $1 + A$, for an array A , denotes adding one to all elements of A . Symmetrically, if $n_{bx} = 0$, then $\text{LCP}[i_{ax} + 1..j_{ax}) = 1 + \text{LCP}[i_x + 1..j_x)$.*

Algorithm 1: Infer BWT from an LCP array

Input: an array $\text{LCP}[1..n)$ of integers and ω 's
Output: a string $v \in \{a, b\}^n$ such that $\text{LCP}_{\text{IBWT}(v)} = \text{LCP}$ together with a set S of swap intervals, or **false** if there is no such string v

```

1  $S := \emptyset$ ;
2 preprocess LCP for RMQs;
3  $k := \text{RMQ}_{\text{LCP}}[1..n)$ ;
4 if  $\text{LCP}[k] \neq 0$  then
5   if  $\text{LCP}[k] = \omega$  then return  $a^n, \emptyset$ ;
6   else return false;
7 InferInterval( $[0, n), [0, k), [k, n)$ );
8 compute  $W = \text{IBWT}(v)$ ,  $\text{SA}_W$ , and  $\text{LCP}_W$ ;
9 if  $\text{LCP}_W \neq \text{LCP}$  then return false;
10 return  $v, S$ ;
```

The main procedure is given in Algorithm 1. The main work is done in the recursive procedure InferInterval given in Algorithm 2. The procedure gets as input the x -, ax - and bx -intervals for some (unknown) string x , splits the x -interval into xya - and xyb -subintervals based on Lemma 3, and tries to split ax - and bx -intervals similarly. If all subintervals are nonempty, the algorithm processes the two subinterval triples recursively (lines 28 and 29).

When trying to split the ax -interval, the result may be, for example, that the $axya$ -interval is empty. In this case, we do not need to recurse on the xya -interval since the corresponding part of v must be all b 's. The algorithm recognizes the emptiness of $axya$ - or $axyb$ -interval by the fact that $m_{ax} > m_x + 1$, but the problem is to decide which is the empty one. In most cases, this can be determined by comparing the sizes of the different subintervals or even the actual LCP-intervals (see Lemma 5).

There is one case, where the algorithm is unable to determine the empty subintervals, which is when $\text{LCP}[i_{ax} + 1..j_{ax}) = \text{LCP}[i_{bx} + 1..j_{bx}) = 1 + \text{LCP}[i_x + 1..k_x) = 1 + \text{LCP}[k_x + 1..j_x)$. Then, either the $axya$ - and $bxyb$ -intervals are empty or the $axyb$ - and $bxya$ -intervals are empty, but there is no way of deciding between the two cases. It turns out that both are valid choices. The algorithm sets v according to one choice (line 8) but records the alternative choice by adding the interval to the set S . In such a case, the string xy is called a *swap core* and the xy -interval (equal to the x -interval) is called a *swap interval*.

For each swap interval $[i..j)$, the algorithm sets $v[i..k) = aa \dots a$ and $v[k..j) = bb \dots b$, where $k = (i + j)/2$, but swapping the two halves would be an equally good choice. Therefore, if the output of the algorithm contains s swap intervals, it represents a set of 2^s distinct strings. The following lemma shows that the swaps indeed do not affect the LCP array (with the proof in the appendix).

Lemma 6. *Let $v \in \{a, b\}^n$, $W = \text{IBWT}(v)$, $\text{SA} = \text{SA}_W$ and $\text{LCP} = \text{LCP}_W$. Let x be a string that occurs in W and satisfies: (1) $\text{LCP}[i_{xa} + 1..j_{xa}) = \text{LCP}[i_{xb} + 1..j_{xb})$, and (2) $v[i_{xa}..j_{xa}) = aa \dots a$ and $v[i_{xb}..j_{xb}) = bb \dots b$, where $[i_z..j_z)$ is the z -interval for $z \in \{xa, xb\}$. Let v' be the same as v except that $v'[i_{xa}..j_{xa}) = bb \dots b$ and $v'[i_{xb}..j_{xb}) = aa \dots a$. Then $\text{LCP}_{\text{IBWT}(v')} = \text{LCP}$.*

Algorithm 2: InferInterval($[i_x..j_x)$, $[i_{ax}..j_{ax})$, $[i_{bx}..j_{bx})$)

Input: (nonempty) x -, ax - and bx -intervals
Output: Set $v[i_x..j_x)$ and add the swap intervals within $[i_x..j_x)$ to S

```
1  $k_x := \text{RMQ}_{\text{LCP}}[i_x + 1..j_x]$ ;  
2  $m_x := \text{LCP}[k_x]$ ;  
3 if  $j_{ax} - i_{ax} = 1$  then  
4    $k_{ax} := i_{ax}$ ;  
5    $m_{ax} := \omega$ ;  
6 else  
7    $k_{ax} := \text{RMQ}_{\text{LCP}}[i_{ax} + 1..j_{ax}]$ ;  
8    $m_{ax} := \text{LCP}[k_{ax}]$ ;  
9 if  $j_{bx} - i_{bx} = 1$  then  
10   $k_{bx} := i_{bx}$ ;  
11   $m_{bx} := \omega$ ;  
12 else  
13   $k_{bx} := \text{RMQ}_{\text{LCP}}[i_{bx} + 1..j_{bx}]$ ;  
14   $m_{bx} := \text{LCP}[k_{bx}]$ ;  
15 if  $m_{ax} > m_x + 1$  and  $m_{bx} > m_x + 1$  then  
16   if  $\text{LCP}[i_{ax} + 1..j_{ax}] = 1 + \text{LCP}[i_x + 1..k_x]$  then  
17      $v[i_x..k_x) = aa \dots a$ ;  
18      $v[k_x..j_x) = bb \dots b$ ;  
19     if  $\text{LCP}[i_{ax} + 1..j_{ax}] = 1 + \text{LCP}[k_x + 1..j_x]$  then  
20        $\text{add } [i_x..j_x) \text{ to } S$ ;  
21   else  
22      $v[i_x..k_x) = bb \dots b$ ;  
23      $v[k_x..j_x) = aa \dots a$ ;  
24 else if  $m_{ax} > m_x + 1$  then  
25   if  $k_{bx} - i_{bx} = k_x - i_x$  then  
26      $v[i_x..k_x) = bb \dots b$ ;  
27     InferInterval( $[k_x..j_x)$ ,  $[i_{ax}..j_{ax})$ ,  $[k_{bx}..j_{bx})$ );  
28   else  
29      $v[k_x..j_x) = bb \dots b$ ;  
30     InferInterval( $[i_x..k_x)$ ,  $[i_{ax}..j_{ax})$ ,  $[i_{bx}..k_{bx})$ );  
31 else if  $m_{bx} > m_x + 1$  then  
32   if  $k_{ax} - i_{ax} = k_x - i_x$  then  
33      $v[i_x..k_x) = aa \dots a$ ;  
34     InferInterval( $[k_x..j_x)$ ,  $[k_{ax}..j_{ax})$ ,  $[i_{bx}..j_{bx})$ );  
35   else  
36      $v[k_x..j_x) = aa \dots a$ ;  
37     InferInterval( $[i_x..k_x)$ ,  $[i_{ax}..k_{ax})$ ,  $[i_{bx}..j_{bx})$ );  
38 else  
39   InferInterval( $[i_x..k_x)$ ,  $[i_{ax}..k_{ax})$ ,  $[i_{bx}..k_{bx})$ );  
40   InferInterval( $[k_x..j_x)$ ,  $[k_{ax}..j_{ax})$ ,  $[k_{bx}..j_{bx})$ );
```

Theorem 1. Algorithm 1 computes in linear time a representation of the set of all strings $v \in \{a, b\}^*$ such that $\text{LCP}_{\text{IBWT}(v)}$ is the input array, or returns false if no such string exists.

Proof. Since the algorithm verifies its result (lines 9 and 10), it will return false if the input is not a valid LCP array. Given a valid LCP array, Algorithm 2 sets all elements of v since it recurses on any subinterval that it doesn't set. All the choices made by the algorithm are forced by the lemmas in this and the previous section. The swap intervals record all alternatives in the

cases where the content of v could not be fully determined, and all of those alternatives have the same LCP array by Lemma 6. It is also easy to see that the algorithm runs in linear time. \square

4 Coupling Constrained Eulerian Cycle

We will now set out to prove the NP-completeness of the single string inference problems BCSILA and BTSILA. The proofs are done by a reduction from 3-SAT via an intermediate problem called Coupling Constrained Eulerian Cycle (CCEC) described in this section.

Consider a directed graph G of degree two, i.e., every vertex in G has exactly two incoming and two outgoing edges. If G is connected, it is Eulerian. An Eulerian cycle can pass through each vertex in two possible ways, which we call the *straight state* and the *crossing state* of the vertex as illustrated here:



We consider each vertex to be a *switch* that can be flipped between these two states. The combination of vertex states is called the *graph state*. For a given graph state, the paths in the graph form, in general, a collection of cycles. The Eulerian cycle problem can then be stated as finding a graph state such that there is only a single cycle; we call such a graph state Eulerian.

In the *Coupling Constrained Eulerian Cycle (CCEC) problem*, we are given a graph as described above, an initial graph state, and a partitioning of the set of vertices. If we flip a vertex state, we must simultaneously flip the states of all the vertices in the same partition, i.e., the vertices in a partition are coupled. A graph state that is achievable from the initial state by a set of such *partition flips* is called a *feasible state*. The CCEC problem is to determine if there exists a feasible graph state that is Eulerian.

Theorem 2. *CCEC is NP-complete.*

Proof. The proof is by reduction from 3-SAT. To obtain a CCEC graph from a 3-CNF formula, a gadget of five vertices is constructed from each clause and these gadgets are connected by a cycle. In each gadget, three of the vertices are labeled by the literals of the corresponding clause; the other two are called free vertices. See Fig. 1 for an illustration.

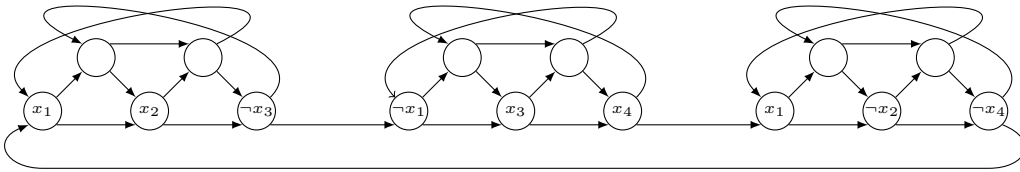


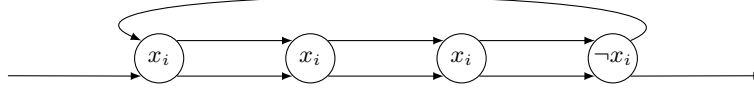
Fig. 1. The CCEC graph corresponding to a 3-CNF formula $(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_3 \vee x_4) \wedge (x_1 \vee \neg x_2 \vee \neg x_4)$.

Each labeled vertex is in a straight state if the labeling literal is false and in a crossing state if the literal is true; their initial state corresponds to some arbitrary truth assignment to the variables. For each variable x_i , there is a vertex partition consisting of all vertices labeled by x_i or $\neg x_i$, so that flipping this partition corresponds to changing the truth value of x_i . Each free vertex forms a singleton partition and has an arbitrary initial state. Thus a graph state is feasible iff the labeled vertex states correspond to some truth assignment.

If a clause is false for a given truth assignment, the labeled vertices in the corresponding gadget are all in a straight state. This separates a part of the gadget from the main cycle and thus the graph state is not Eulerian. If a clause is true, at least one of the labeled vertices in the

gadget is in a crossing state. Then we can always choose the state of the free vertices so that the full gadget is connected to the main cycle. Thus there exists a feasible Eulerian graph state iff there exists a truth assignment to the variables that satisfies all clauses. \square

For purposes that will become clear later, we modify the above construction by adding some extra components to the graph without changing the validity of the reduction. Specifically, for each variable x_i in the 3-CNF formula we add the following gadget to the main cycle:



The vertices in the gadget are treated similarly to the other vertices in the graph: they belong to the partition with the other vertices labeled by x_i or $\neg x_i$, and the initial state is determined by the truth value of the labeling literal. It is easy to see that the gadget will be fully connected to the main cycle whether x_i is true or false. Thus the extra gadgets have no effect on the existence of an Eulerian cycle. Finally, we insert to the main cycle a single vertex labelled y with a self loop and forming a singleton partition.

5 BCSILA to CCEC

The next step is to establish a connection between the BCSILA and CCEC problems by showing a reduction from BCSILA to CCEC. Although the direction of the reduction is opposite to what we want, this construction plays a key role in the analysis of the main construction described in the next section.

Given a BCSILA instance (an integer array), we use Algorithm 1 to produce a representation of a set V of strings. The problem is then to decide if there exists $v \in V$ such that $\text{IBWT}(v)$ is a single (cyclic) string. We will write V as a string with brackets marking the swaps. For example, $V = b[ab][ab]a = \{bababa, babbaa, bbaaba, bbabaa\}$. In Example 1, we saw that the inverse BWT of a string $v \in V$ can be represented as a graph G_v where the vertices are labeled by positions in v and there is an edge between vertices i and j if, for some character $c \in \{a, b\}$ and some integer k , $\hat{v}[i] = c$ is the k th occurrence of c in \hat{v} and $v[j] = c$ is the k th occurrence of c in v . Such an edge (i, j) is labeled by c_k . Note that $\forall v \in V$, \hat{v} is the same; we will denote it by \hat{V} . We form a generalized graph G_V as a union of the graphs G_v , $v \in V$ (see Fig. 11 for an example).

Consider a_k (the k th a) in \hat{V} , say at position i . If a_k is outside any swap region in V , say at position j , there is a single edge (i, j) in G_V labeled by a_k . If a_k is within a swap region in V , it has two possible positions in the strings $v \in V$, say j and j' . That same pair of positions are also the possible positions of some b , say $b_{k'} = \hat{V}[i']$. Then G_V has two edges, (i, j) and (i, j') , labeled with a_k and two edges, (i', j) and (i', j') , labeled with $b_{k'}$. The positions/vertices j and j' are called a *swap pair*.

To obtain a CCEC graph \tilde{G}_V , we make two modifications to G_V . First, we merge each swap pair into a single vertex. Each merged vertex now has two incoming and two outgoing edges and all other vertices have one incoming and one outgoing edge. Second, we remove all vertices with degree one by concatenating their incoming and outgoing edges (see Fig. 11).

The initial state of the vertices in \tilde{G}_V is set so that the cycles in \tilde{G}_V correspond to the cycles in G_v for some $v \in V$. Two vertices in \tilde{G}_V belong to the same partition if their labels belong to the same swap interval in V . Then we have a one-to-one correspondence between swaps in V and partition flips in \tilde{G}_V . If this CCEC instance has a solution, the Eulerian cycle spells a single string realizing the input LCP array. If the CCEC instance has no solution, the original BCSILA problem has no solution either.

6 BCSILA is NP-Complete

We are now ready to show that BCSILA is NP-complete using the reduction chain $3\text{-SAT} \rightarrow \text{CCEC} \rightarrow \text{BCSILA}$. The first step was described in Section 4, and we will next describe the second. The latter reduction is not a general reduction from an arbitrary CCEC instance but works only for a CCEC instance obtained by the first reduction (including the extra gadgets).

The above BCSILA to CCEC reduction transforms each pair of swapped positions into a vertex and each swap interval into a vertex partition. Our construction creates a BCSILA instance such that the resulting BWT has the necessary swaps to produce the CCEC instance vertices and partitions. However, the BWT also has some unwanted swaps producing spurious vertices, but we will show that these spurious vertices do not invalidate the reduction.

Starting from a CCEC instance, we construct a set of cyclic strings and obtain the BCSILA instance as the LCP array of that string set. The construction associates two strings to each vertex and the cyclic strings are formed by concatenating the vertex strings according to the cycles in the graph in its initial state. The two passes of the cycles through a vertex must use different strings but it does not matter which pass uses which string.

Let n be the number of vertices in the CCEC graph and let m be the number of vertex partitions. We number the vertices from 1 to n and the partitions from 1 to m . The biggest partition number is assigned to the partition with the vertex y , the second biggest to the partition corresponding to the variable x_1 , the third biggest to variable x_2 , and so on. The three biggest vertex numbers are assigned to the vertices labeled x_1 in the extra gadget for the variable x_1 , the next three biggest to the extra gadget vertices labeled x_2 and so on. Within each extra gadget, the biggest number is assigned to the middle one of the three vertices. The strings associated with a vertex are ba^kba^{m+2h} and bba^kbba^{m+2h-1} , where k is the partition number and h is the vertex number. This completes the description of the transformation from a CCEC instance to a BCSILA instance.

Let us now analyze the transformation by changing the BCSILA instance back to a CCEC instance using the construction of the preceding section. Specifically, we will analyze the swaps in the BWT produced from the LCP array. Let W be the set of cyclic strings constructed from the CCEC instance, and let V be the BWT with swaps constructed from LCP_W . An interval $[i..j]$ in V is a swap interval if and only if (1) $[i..j]$ is an x -interval for a string x such that either $\text{occ}(axa) = \text{occ}(bxb) = \text{occ}(x)/2$ or $\text{occ}(axb) = \text{occ}(bxa) = \text{occ}(x)/2$, where $\text{occ}(y)$ is the number of occurrences of y in W , and (2) $\text{LCP}_W[i+1..k] = \text{LCP}_W[k+1..j]$, where $k = (i+j)/2$. If $[i..j]$ is a swap interval, the string x is called its *swap core*. Our goal is to identify all swap cores.

Let us first consider strings of the form $x = ba^kb$. If $k > m$, $\text{occ}(x) \leq 1$ and x cannot be a swap core. For $k \in [1..m]$, x is always a swap core and corresponds to the CCEC partition numbered k . Let $v = \text{BWT}(W)$ and let V' be v together with the swaps for cores of the form $x = ba^kb$, $k \in [1..m]$. It is easy to verify that a CCEC instance constructed from V' as described in the previous section is identical to the original CCEC instance. Thus, if there were no other swap cores, we would have a perfect reduction.

Unfortunately, there are other swap cores. A systematic examination of all strings in Appendix F shows that the other swap cores must be of the following forms: ba^{m+2n-1} , $a^{m+2n-1}b$, a^mba^m , a^mbba^m , a^kba^h , a^kbba^h , $a^kba^iba^h$ and $a^kbba^ibba^h$. Furthermore, it shows that each such swap core has exactly two occurrences, which means that the values k and/or h have to be sufficiently large. Each extra swap core adds a free vertex that is connected to the graph by making two existing edges to pass through the new vertex. Because of the way we chose to assign the biggest partition and vertex numbers, all the additional connections are within the extra gadgets, which does not change the existence of an Eulerian cycle. This completes the proof.

Theorem 3. *BCSILA is NP-complete.*

7 BTSILA is NP-Complete

We will now show that BTSILA is NP-complete by modifying the above reduction for BCSILA to include a single terminator symbol $\$$ in the strings. The modification is applied to the set W of cyclic strings derived from the CCEC instance such that LCP_W is the BCSILA instance. Specifically, we replace the (unique) occurrence of a^{m+2n} , which is the longest consecutive run of a 's, with $a^{m+2n+1}\$a^{m+2n}$ to obtain $W_\$$ and $\text{LCP}_{W_\$}$. We will show that $\text{LCP}_{W_\$}$ is a yes-instance of CSILA iff LCP_W is a yes-instance of BCSILA. Furthermore, if a cyclic string u is a solution to the CSILA instance, i.e., $\text{LCP}_u = \text{LCP}_{W_\$}$, then $\text{LCP}_v = \text{LCP}_{W_\$}$, where v is the rotation of u ending with $\$$ interpreted as a terminated string. Thus $\text{LCP}_{W_\$}$ is a yes-instance of BTSILA iff it is a yes-instance of CSILA iff LCP_W is a yes-instance of BCSILA.

In general, adding even a single occurrence of a third symbol complicates the inference of the BWT from the LCP array and means that the set of equivalent BWTs can no more be described by a set of swaps. Consider how the operation of the procedure `InferInterval` (Algorithm 2) changes. First, it gets an extra $\$x$ -interval as an input in addition to x -, ax - and bx -intervals. Second, the x -interval may be split into three subintervals, $xy\$$ -, xya - and xyb -intervals, instead of two (which happens when the LCP interval contains two identical minima). This leads to many more combinations to consider, and some of those combinations are more complicated.

Fortunately, in our case, having the single $\$$ surrounded by the two longest runs of a 's simplifies things, and we will describe a modification of `InferInterval` to handle this case. Every call to `InferInterval` belongs to one of the following three types: (1) the x -interval is split into two and the $\$x$ -interval is empty, (2) the x -interval is split into two and the $\$x$ -interval is non-empty, and (3) the x -interval is split into three. The first case needs no modification at all. The other two cases mean that either $\$x$ or $x\$$ occurs in the produced string set, and since this property is not affected by swaps (or the threeway permutations described below), one of them occurs in every produced string set including $W_\$$. Since x must occur at least twice, one of the latter two cases happens iff $x = a^k$ for some $k \in [0..m+2n]$. Although in general `InferInterval` cannot always know x , it is easy to keep track of x when $x = a^k$.

When `InferInterval` is called with $x = a^k$ for $k \leq m+2n-2$, the x -interval and the ax -interval are always split into three, the bx -interval is split into two, and there is a $\$x$ -interval of size one. In general, we might not know whether the two subintervals of bx -interval are $bx\$$ - and bxa -, or $bx\$$ - and $bx b$ -, or bxa - and $bx b$ -intervals. However, since $x\$$ - and $ax\$$ -intervals both have size one, there can be no $bx\$$ -interval, and thus all the subintervals can be uniquely determined and recursed on. When $x = a^{m+2n-1}$, the x -interval has size five and is split into three with the middle part (xa -interval) having size three. The ax interval has size three and is split into three. In this case too, only one combination of subintervals is possible.

When $x = a^{m+2n}$, the x -interval has size three and is split into three, and the $\$x$ -, ax - and bx -intervals have size one. Therefore, the x -interval in the BWT contains some permutation of the three characters and all permutations are valid. This threeway permutation adds to the variation provided by the swaps in other parts of the BWT. A more careful analysis shows that the BWT x -interval of

- $\$ab$ or $\$ba$ implies an occurrence of $\$x\$$ which is only possible if $x\$$ is a separate string;
- $ba\$$ implies an occurrence of axa which is only possible if a single a is separate string;
- $a\$b$ implies occurrences of $ax\$$ and $\$xa$ which is only possible if $ax\$$ is a separate string;
- $ab\$$ implies an occurrence of $ax\$xb$; and
- $b\$a$ implies an occurrence of $bx\$xa$.

A single string solution is only possible in the last two cases, and any such solution corresponds to a solution for the BCSILA instance LCP_W (obtained by replacing $ax\$x$ or $x\$ax$ with x). Hence

$\text{LCP}_{W_\$}$ is a yes-instance of CSILA, and thus of BTSILA, if and only if LCP_W is a yes-instance of BCSSILA, which proves the following result.

Theorem 4. *BTSILA is NP-complete.*

8 Algorithm for CSSILA

In all of the above, we have assumed a binary alphabet (excluding the single symbol \$). In this section, we consider the CSSILA problem (i.e. Cyclic String Set Inference from LCP Array) without a restriction on the alphabet size.

Let $\mathcal{L}[1..n]$ be an instance of the CSSILA problem, i.e., an array of integers (and possibly ω 's). Let $\sigma - 1$ be the number of zeroes in \mathcal{L} , and Σ an alphabet of size σ . As with the binary BCSSILA problem, we describe an algorithm that outputs a representation of the set $W_{\mathcal{L}} = \{w \in \Sigma^n : \text{LCP}_{\text{IBWT}(w)} = \mathcal{L}\}$; in this case the representation is an automaton that accepts $W_{\mathcal{L}}$. We show the following result.

Theorem 5. *Given an array $\mathcal{L}[1..n]$ of integers (and possibly ω 's) containing $\sigma - 1$ zeroes, we can construct a deterministic finite automaton recognizing $W_{\mathcal{L}}$ in time $O(\sigma^2 2^\sigma (\frac{n}{\sigma} + 1)^\sigma)$ and space $O(\sigma 2^\sigma (\frac{n}{\sigma} + 1)^\sigma)$.*

The algorithm and further details are in Appendix D.

References

1. Mohamed Ibrahim Abouelhoda, Stefan Kurtz, and Enno Ohlebusch. Replacing suffix trees with enhanced suffix arrays. *Journal of Discrete Algorithms*, 2(1):53–86, 2004.
2. Amihood Amir. Personal communication, String Masters in Rouen, France, 3–5 February, 2014.
3. Alberto Apostolico. The myriad virtues of subword trees. In Alberto Apostolico and Zvi Galil, editors, *Combinatorial Algorithms on Words*, NATO ASI Series F12, pages 85–96. Springer-Verlag, Berlin, Germany, 1985.
4. Alberto Apostolico, Maxime Crochemore, Martin Farach-Colton, Zvi Galil, and S. Muthukrishnan. 40 years of suffix trees. *Commun. ACM*, 59(4):66–73, 2016.
5. Hideo Bannai, Shunsuke Inenaga, Ayumi Shinohara, and Masayuki Takeda. Inferring strings from graphs and arrays. In Branislav Rovan and Peter Vojtás, editors, *Mathematical Foundations of Computer Science 2003, 28th International Symposium, MFCS 2003, Bratislava, Slovakia, August 25–29, 2003, Proceedings*, volume 2747 of *Lecture Notes in Computer Science*, pages 208–217. Springer, 2003.
6. Bastien Cazaux and Eric Rivals. Reverse engineering of compact suffix trees and links: A novel algorithm. *J. Discrete Algorithms*, 28:9–22, 2014.
7. Julien Clément, Maxime Crochemore, and Giuseppina Rindone. Reverse engineering prefix tables. In Susanne Albers and Jean-Yves Marion, editors, *26th International Symposium on Theoretical Aspects of Computer Science, STACS 2009, February 26–28, 2009, Freiburg, Germany, Proceedings*, volume 3 of *LIPIcs*, pages 289–300. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, 2009.
8. Maxime Crochemore and Lucian Ilie. Computing longest previous factor in linear time and applications. *Inf. Process. Lett.*, 106(2):75–80, 2008.
9. Maxime Crochemore, Costas S. Iliopoulos, Solon P. Pissis, and German Tischler. Cover array string reconstruction. In Amihood Amir and Laxmi Parida, editors, *Combinatorial Pattern Matching, 21st Annual Symposium, CPM 2010, New York, NY, USA, June 21–23, 2010. Proceedings*, volume 6129 of *Lecture Notes in Computer Science*, pages 251–259. Springer, 2010.
10. Jean-Pierre Duval, Thierry Lecroq, and Arnaud Lefebvre. Border array on bounded alphabet. *Journal of Automata, Languages and Combinatorics*, 10(1):51–60, 2005.
11. Jean-Pierre Duval, Thierry Lecroq, and Arnaud Lefebvre. Efficient validation and construction of border arrays and validation of string matching automata. *RAIRO-Theor. Inf. Appl.*, 43(2):281–297, 2009.
12. František Franěk, S. Gao, Weilin Lu, Patrick J. Ryan, William F. Smyth, Yu Sun, and Lu Yang. Verifying a border array in linear time. *Journal on Combinatorial Mathematics and Combinatorial Computing*, 42:223–236, 2002.

13. Pawel Gawrychowski, Artur Jez, and Lukasz Jez. Validating the knuth-morris-pratt failure function, fast and online. *Theory Comput. Syst.*, 54(2):337–372, 2014.
14. Ira M. Gessel and Christophe Reutenauer. Counting permutations with given cycle structure and descent set. *Journal of Combinatorial Theory, Series A*, 64(2):189–215, 1993.
15. Dan Gusfield. *Algorithms on Strings, Trees, and Sequences : Computer Science and Computational Biology*. Cambridge University Press, Cambridge, United Kingdom, 1997.
16. Jing He, Hongyu Liang, and Guang Yang. Reversing longest previous factor tables is hard. In Frank Dehne, John Iacono, and Jörg-Rüdiger Sack, editors, *Algorithms and Data Structures - 12th International Symposium, WADS 2011, New York, NY, USA, August 15-17, 2011. Proceedings*, volume 6844 of *Lecture Notes in Computer Science*, pages 488–499. Springer, 2011.
17. Peter M. Higgins. Burrows-Wheeler transformations and de Bruijn words. *Theor. Comput. Sci.*, 457:128–136, 2012.
18. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Verifying and enumerating parameterized border arrays. *Theor. Comput. Sci.*, 412(50):6959–6981, 2011.
19. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Inferring strings from suffix trees and links on a binary alphabet. *Discrete Applied Mathematics*, 163:316–325, 2014.
20. Juha Kärkkäinen, Dominik Kempa, and Marcin Piątkowski. Tighter bounds for the sum of irreducible LCP values. *Theoretical Computer Science*, 2015.
21. Gregory Kucherov, Lilla Tóthmérés, and Stéphane Viallette. On the combinatorics of suffix arrays. *Inf. Process. Lett.*, 113(22-24):915–920, 2013.
22. Udi Manber and Gene W. Myers. Suffix arrays: a new method for on-line string searches. *SIAM J. Comp.*, 22(5):935–948, 1993.
23. Sabrina Mantaci, Antonio Restivo, Giovanna Rosone, and Marinella Sciortino. An extension of the Burrows-Wheeler transform. *Theor. Comput. Sci.*, 387(3):298–312, 2007.
24. Yuto Nakashima, Takashi Okabe, Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Inferring strings from Lyndon factorization. In Erzsébet Csuhaj-Varjú, Martin Dietzfelbinger, and Zoltán Ésik, editors, *Mathematical Foundations of Computer Science 2014 - 39th International Symposium, MFCS 2014, Budapest, Hungary, August 25-29, 2014. Proceedings, Part II*, volume 8635 of *Lecture Notes in Computer Science*, pages 565–576. Springer, 2014.
25. Enno Ohlebusch. *Bioinformatics Algorithms: Sequence Analysis, Genome Rearrangements, and Phylogenetic Reconstruction*. Oldenbusch Verlag, 2013.
26. Nicolas Philippe. Caractérisation et énumération des arbres compacts des suffixes. Master’s thesis, Université de Rouen, 2007.
27. Klaus-Bernd Schürmann and Jens Stoye. Counting suffix arrays and strings. *Theor. Comput. Sci.*, 395(2-3):220–234, 2008.
28. Imre Simon. Piecewise testable events. In *Proc. 2nd GI Conference on Automata Theory and Formal Languages*, volume 33 of *LNCS*, pages 214–222. Springer, 1975.
29. Bill Smyth. *Computing Patterns in Strings*. Pearson Addison-Wesley, Essex, England, 2003.
30. Tatiana A. Starikovskaya and Hjalte Wedel Vildhøj. A suffix tree or not a suffix tree? *J. Discrete Algorithms*, 32:14–23, 2015.
31. Peter Weiner. Linear pattern matching algorithms. In *Proceedings of the 14th Annual IEEE Symposium on Switching and Automata Theory*, pages 1–11, 1973.

A Examples of Suffix and LCP Arrays and Suffix Trees

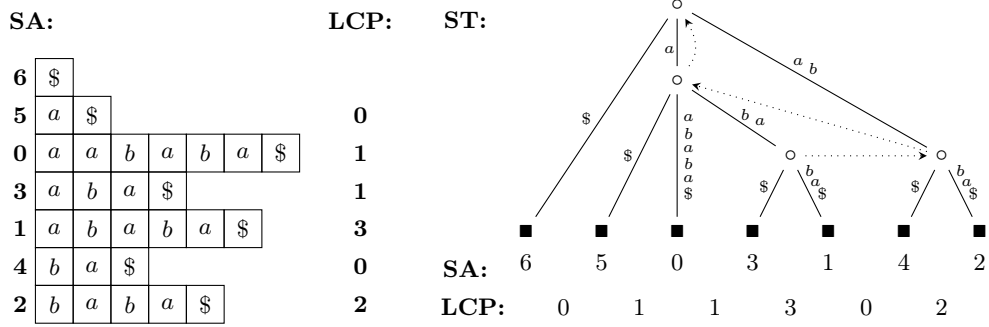


Fig. 2. SA, LCP and ST for terminated string *aababa\$*. Notice how the LCP array encodes the shape of the suffix tree. The dashed arrows are suffix links, which connect node representing *cx* for a symbols *c* and a string *x* to node representing *x*.

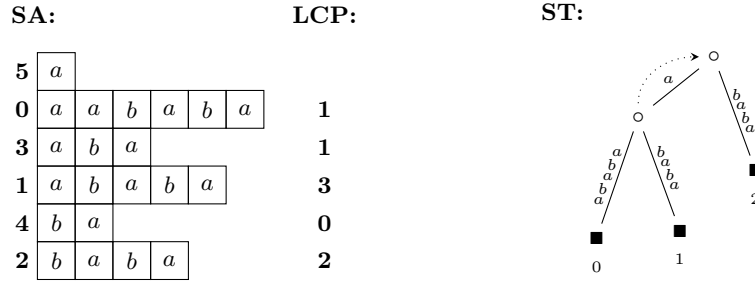


Fig. 3. SA, LCP and ST for open-ended string *aababa*. In the suffix tree, the suffixes that are proper prefixes of another suffixes are not represented by a leaf (or another special node).

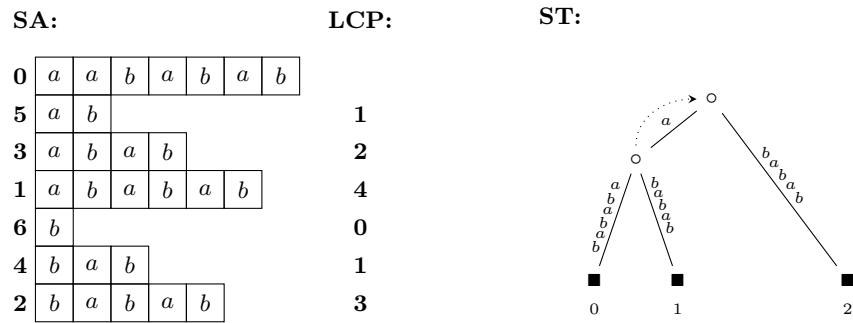


Fig. 4. SA, LCP and ST for open-ended string *aababab*, which is an extension of the string in the preceding figure. The suffix tree shape is the same but SA and LCP are quite different.

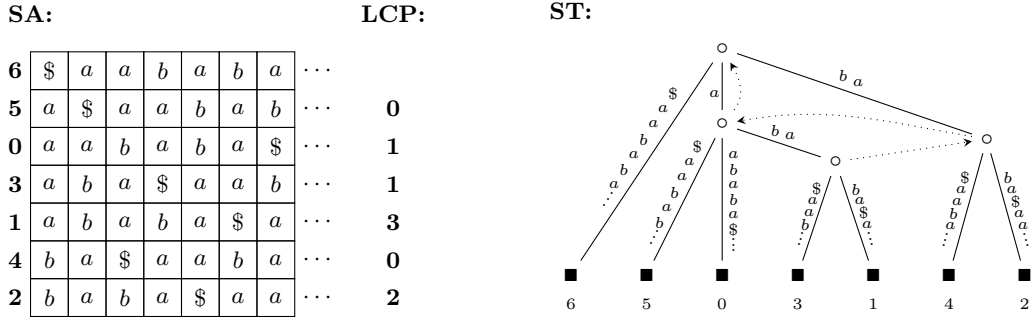


Fig. 5. SA, LCP and ST for cyclic string *aababa*\$. Notice that SA, LCP and suffix tree shape are the same as in Fig. 2

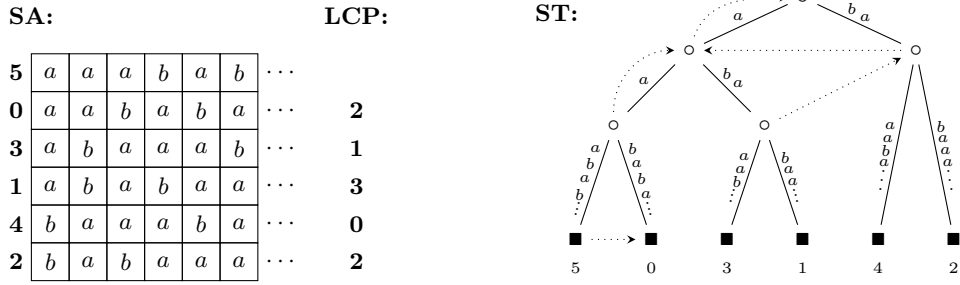


Fig. 6. SA, LCP and ST for cyclic string *aababa*.

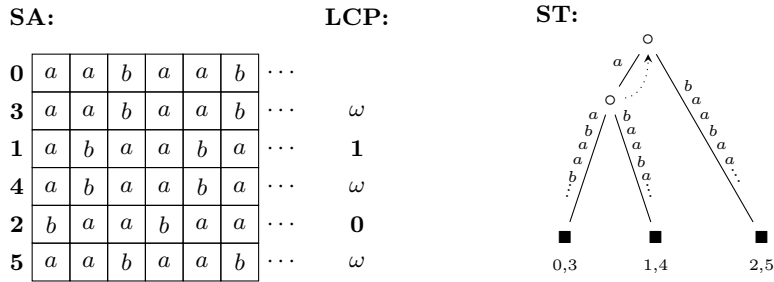


Fig. 7. SA, LCP and ST for cyclic string *aabaab*. Because the string is non-primitive (concatenation of multiple copies of the same string), some of its cyclic suffixes are identical. The LCP of identical suffixes is ω and they share a leaf in the suffix tree.

B Reductions from BTSILA

Proof (of Proposition 1). By the discussion in the introduction, an array of n integers is

- a yes-instance of BTSILA iff it has a leading zero and is a yes-instance of BOSILA with the leading zero removed,
- a yes-instance of BTSILA iff it has a leading zero and at most one other zero, and is a yes-instance of TSILA,
- a yes-instance of TSILA iff it has a leading zero and is a yes-instance of OSILA with the leading zero removed,
- a yes-instance of TSILA iff it is a yes-instance of TSSILA,
- a yes-instance of TSSILA iff it has one or more leading zeros and is a yes-instance of OSSILA with the leading zeros removed,
- a yes-instance of BTSILA iff it has a leading zero and at most one other zero, and is a yes-instance of BTSSILA, and
- a yes-instance of BTSSILA iff it has one or more leading zeros and at most one other zero, and is a yes-instance of BOSSILA with the leading zeros removed.

In all cases, there is a simple linear or at most quadratic time reduction. \square

C Algorithm for BCSSILA: A Proof and an Example

Proof (of Lemma 6). Consider first how $\Psi_{v'}$ differs from Ψ_v . For any $i \in [0..n)$, if $\Psi_v[i] \notin [i_x..j_x)$ then $\Psi_{v'}[i] = \Psi_v[i]$. Otherwise $\Psi_{v'}[i] = \Psi_v[i] + n_{xa} \in [i_x..j_x)$ or $\Psi_{v'}[i] = \Psi_v[i] - n_{xa} \in [i_x..j_x)$, i.e., it is swapped from one side of the interval $[i_x..j_x)$ to the other side.

Now we use Lemma 1 to determine how a suffix at $\text{SA}[i]$ changes with the swap. If i belongs to a cycle that never visits $[i_x..j_x)$, i.e., the suffix does not contain x , there is no change. Suppose then that the cycle starting at i first reaches $[i_x..j_x)$ after k steps, and w.l.o.g. assume that it reaches specifically the xa -interval, i.e. $\Psi_v^k[i] \in [i_{xa}..j_{xa})$. Then for some string y of length k , the suffix at i changes from $yxa\dots$ into $yxb\dots$. Note also that yx cannot contain x except at the end.

Now consider two adjacent suffixes. If both are of the form $yxa\dots$, they both change to $yxb\dots$. The parts after x may change a lot but LCP of the two suffixes remains the same because $\text{LCP}[i_{xa} + 1..j_{xa}) = \text{LCP}[i_{xb} + 1..j_{xb})$. In all other cases (one or both do not contain x or the parts before x differ), the LCP is determined in the unchanged part of the suffixes. Thus $\text{LCP}_{\text{IBWT}(v')} = \text{LCP}$. \square

The following example illustrates the operation of the algorithm.

Example 3. Let us consider an integer array $\mathcal{L}[1..7) = [1, 4, 0, 2, 1, 3]$. Using the above algorithms we will try to reconstruct a string v , such that $\text{LCP}_{\text{IBWT}(v)} = \mathcal{L}$. Since $\mathcal{L}[3] = 0$ w contains 3 occurrences of a and 4 occurrences of b , and the initial call to Algorithm 2 is $\text{InferInterval}([0..7), [0..3), [3..7))$ (see Figure 8 (1)). We then have $m_x = \mathcal{L}[3] = 0$, $m_{ax} = \mathcal{L}[1] = 1$ and $m_{bx} = \mathcal{L}[5] = 1$, which leads to the recursive calls $\text{InferInterval}([0..3), [0..1), [3..5))$ and $\text{InferInterval}([3..7), [1..3), [5..7))$.

When processing $\text{InferInterval}([0..3), [0..1), [3..5))$ (see Figure 8 (2)), we find that $m_{bx} = m_x + 1 = 2$ but $m_{ax} = \omega$ because the ax -interval has size 1. Thus we set $v[0..1) = b$ (line 14) and make the recursive call $\text{InferInterval}([1..3), [0..1), [4..5))$.

When processing $\text{InferInterval}([1..3), [0..1), [4..5))$ (see Figure 8 (3)), we find that both the ax - and the bx -interval have size 1. In such a case, we always have a swap interval. Here we set $v[1..3) = ab$ and add $[1..3)$ into S .

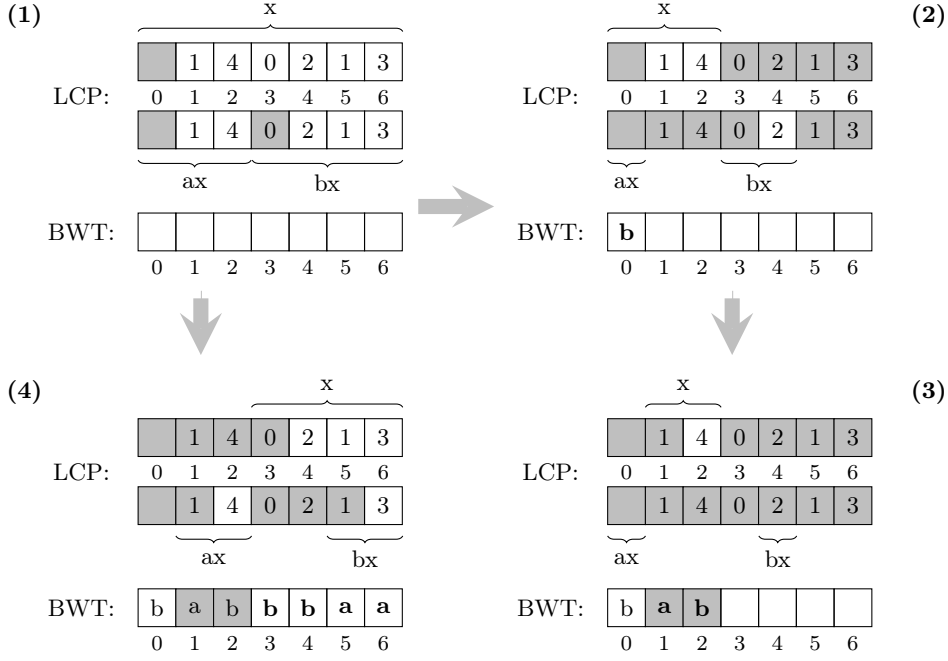


Fig. 8. Graphical illustration of Example 3.

When processing $\text{InferInterval}([3..7], [1..3], [5..7])$ (see Figure 8 (4)), we have $m_x = \mathcal{L}[5] = 1$ but $m_{ax} = 4 > m_x + 1$ and $m_{bx} = 3 > m_x + 1$. Comparing $\mathcal{L}[2..3] = [4]$ and $\mathcal{L}[4..5] = [2]$ (line 10), we find that they do not match. Thus we set $v[3..5] = bb$ and $v[5..7] = aa$.

The final result is $v = b[ab]bbaa$, where the only swap interval is marked with brackets. The main algorithm then computes $W = \text{IBWT}(v) = \{aabb, abb\}$, verifies that $\text{LCP}_W = \mathcal{L}$ and outputs $b[ab]bbaa$. It is easy to verify that $\text{LCP}_{\text{IBWT}(bbabbaa)} = \mathcal{L}$ too.

D Algorithm for CSSILA

In this section we present the algorithm solving CSSILA problem for alphabets of any size. Let $\Sigma = \{a_1, a_2, \dots, a_\sigma\}$ be an alphabet and $\mathcal{L}[1..n]$ be LCP array containing $\sigma - 1$ zeroes. We try to reconstruct a set of strings $W_{\mathcal{L}} = \{w \in \Sigma^n : \text{LCP}_{\text{IBWT}(w)} = \mathcal{L}\}$. The resulting set $W_{\mathcal{L}}$ is represented as an acyclic deterministic finite automaton $\mathcal{A}_{\mathcal{L}}$ accepting all strings $w \in W_{\mathcal{L}}$. Such a representation allows us to perform efficient $W_{\mathcal{L}}$ membership tests, enumerate all its members, and efficiently find the lexicographic predecessor and successor for any $w \in W_{\mathcal{L}}$.

The resursive iteration of intervals in the binary case does not work for larger alphabets, because we can no more uniquely match intervals. Instead, the algorithm iterates from left to right, and for that we need a different characterization of $W_{\mathcal{L}}$.

For any $c \in \Sigma$ and any $w \in W_{\mathcal{L}}$, consider two consecutive occurrences of c in w (i.e., there are no other occurrences of c between them but there may be other characters). Say, they occur at positions h and k , and are the i^{th} and $(i+1)^{\text{th}}$ occurrence of c in w . Then we must have that

$$\mathcal{L}[i_c + i] = 1 + \min\{\mathcal{L}[j] : h < j \leq k\}$$

where i_c is the starting position of the c -interval. We call this the *pair constraint*. The following lemma shows how to characterize $W_{\mathcal{L}}$ using pair constraints.

Lemma 7. *For any $w \in \Sigma^n$, $w \in W_{\mathcal{L}}$ if and only if every pair of consecutive occurrences satisfies the pair constraint.*

Proof. Let $V = \text{IBWT}(w)$. Consider a pair of consecutive occurrences at positions h and k in w , which are the i^{th} and $(i+1)^{\text{th}}$ occurrence of c in w . Let $x = V_{\text{SA}_V[h]}$ and $y = V_{\text{SA}_V[k]}$. Then we must have that $cx = V_{\text{SA}_V[i_c+i-1]}$ and $cy = V_{\text{SA}_V[i_c+i]}$, where i_c is the starting position of the c -interval.

For any suffix array SA and the corresponding LCP array LCP , and any two positions h and k with $h < k$, $\min\{\text{LCP}[j] : h < j \leq k\}$ is the length of the longest common prefix of the suffixes $\text{SA}[h]$ and $\text{SA}[k]$. Thus if $\mathcal{L} = \text{LCP}_V$, we must have

$$\mathcal{L}[i_c + i] = \text{lcp}(cx, cy) = 1 + \text{lcp}(x, y) = 1 + \min\{\mathcal{L}[j] : h < j \leq k\}.$$

This proves the “only if” part.

The “if” part is proven by contradiction. Suppose that all the pair constraints hold in \mathcal{L} but $\mathcal{L} \neq \text{LCP}_V$. Let $\mathcal{L}[d] \neq \text{LCP}_V[d]$ be the smallest wrong value in \mathcal{L} . Assume $\mathcal{L}[d] < \text{LCP}_V[d]$; otherwise we swap the roles of \mathcal{L} and LCP_V and pick the smallest value in LCP_V that differs from \mathcal{L} . Let $c \in \Sigma$ be the character such that d is in the c -interval $[i_c, j_c)$, and let $i = d - i_c$. Let h and k be the positions of the i^{th} and $(i+1)^{\text{th}}$ occurrences of c in w . Since the pair constraints hold for both \mathcal{L} and LCP_V , we must have

$$\begin{aligned} \mathcal{L}[d] &= 1 + \min\{\mathcal{L}[j] : h < j \leq k\} \text{ and} \\ \text{LCP}_V[d] &= 1 + \min\{\text{LCP}_V[j] : h < j \leq k\}. \end{aligned}$$

Let $j \in [h+1..k]$ be a position where $\mathcal{L}[j]$ is minimized in that range, i.e., $\mathcal{L}[j] = \mathcal{L}[d] - 1$. But then we must have $\mathcal{L}[j] < \text{LCP}_V[j]$, which contradicts $\mathcal{L}[d]$ being the smallest wrong value. This completes the “if” part. \square

Recall that for a string $w \in \Sigma^*$ and $c \in \Sigma$, $|w|_c$ denotes the number of occurrences of c in w . We extend this notions to LCP arrays. Namely, $|\mathcal{L}|_c$ denotes the number of occurrences of c in any string w such that $\text{LCP}_w = \mathcal{L}$. We split LCP array \mathcal{L} into σ so-called *character arrays* as follows. For any $c \in \Sigma$, let $[i_c, j_c)$ be the c -interval and let $\mathcal{L}_c[1..j_c - i_c] = \mathcal{L}[i_c + 1..j_c] - 1$ (where $A - 1$ means subtracting one from each element of A). Notice that the c -intervals can be determined solely based on the occurrences of zeroes in \mathcal{L} , and thus we can extend the above definitions to cases where \mathcal{L} is not a valid LCP array. For a technical reason, to avoid a number of special cases to be checked (e.g. for empty character subsequences or boundary cases), we set $\mathcal{L}[0] = \mathcal{L}_c[0] = -1$ and $\mathcal{L}[n] = \mathcal{L}_c[|\mathcal{L}|_c] = -2$ for all $c \in \Sigma$. This gives us a trivial match for the begin and end of each character sequence with the global sequence \mathcal{L} .

To be able to construct the set $W_{\mathcal{L}}$ iteratively we define a notion of (*prefix*) *consistency* of a string $s \in \Sigma^k$ ($k \leq n$) with an LCP array \mathcal{L} when s is considered to be a prefix of some string in $W_{\mathcal{L}}$. For any $c \in \Sigma$, let $\ell_c(s) = \max\{j < k : s[j] = c\} \cup \{-1\}$ be the position of the last occurrence of c in s (or -1 if $|s|_c = 0$). For any $c \in \Sigma$ such that $|s|_c < |\mathcal{L}|_c$, a *partial pair constraint* is

$$\mathcal{L}[i_c + |s|_c] \leq 1 + \min\{\mathcal{L}[j] : \ell_c(s) < j \leq k\}.$$

In other words, it is a pair condition on the pair consisting of the last occurrence of c in s and the next occurrence of c after the end of s . Since we do not know the location of the next occurrence, we only verify that nothing in $\mathcal{L}[0..k]$ violates the condition. Therefore, we have the inequality in place of the equality in the condition.

Definition 3. Let $s \in \Sigma^k$ for $k \leq n$. We say that s is *prefix consistent* with \mathcal{L} , if

1. the pair constraint holds for every pair of consecutive occurrences in s , and
2. the partial pair constraint holds for each $c \in \Sigma$ such that $|s|_c < |\mathcal{L}|_c$.

From the definition and Lemma 7, we immediately get the following.

Corollary 2. $w \in W_{\mathcal{L}}$ if and only if $|w| = n$ and w is prefix consistent with \mathcal{L} .

See Examples 4 and 5 for illustration of strings consistent and inconsistent with a given LCP array.

Let $p(s) = (|s|_{a_1}, |s|_{a_2}, \dots, |s|_{a_\sigma})$ be the Parikh vector of s and $p(\mathcal{L}) = (|\mathcal{L}|_{a_1}, |\mathcal{L}|_{a_2}, \dots, |\mathcal{L}|_{a_\sigma})$ be the Parikh vector of \mathcal{L} . Define

$$b_c(s) = \begin{cases} -1 & \text{if } \mathcal{L}_c[|s|_c] > \min\{\mathcal{L}[j] : \ell_c(s) < j \leq |s|\} \\ 0 & \text{if } \mathcal{L}_c[|s|_c] = \min\{\mathcal{L}[j] : \ell_c(s) < j \leq |s|\} \\ 1 & \text{if } \mathcal{L}_c[|s|_c] < \min\{\mathcal{L}[j] : \ell_c(s) < j \leq |s|\} \end{cases}$$

and $b(s) = (b_{a_1}(s), b_{a_2}(s), \dots, b_{a_\sigma}(s))$. The following is easy to verify.

Lemma 8. A string s violates a partial pair constraint if and only if $b(s)$ contains -1 .

The significance of the vectors $p(s)$ and $b(s)$ is shown by the following lemma.

Lemma 9. Let $s \in \Sigma^k$, $k < n$, be a string prefix consistent with \mathcal{L} . Given $p(s)$ and $b(s)$ (but not s), and $c \in \Sigma$, we can determine whether sc is prefix consistent with \mathcal{L} and compute $p(sc)$ and $b(sc)$ in $O(\sigma)$ time.

Proof. Let us first look at updating the vectors. Let $s \in \Sigma^k$ be a string consistent with an LCP array \mathcal{L} and $a_i \in \Sigma$. Given $p(s) = (|s|_{a_1}, |s|_{a_2}, \dots, |s|_{a_\sigma})$ we have $p(s \cdot a_i) = (|s|_{a_1}, \dots, |s|_{a_{i-1}}, |s|_{a_i} + 1, |s|_{a_{i+1}}, \dots, |s|_{a_\sigma})$.

By definition of b and consistency of s with \mathcal{L} we have:

$$b_{a_i}(s \cdot a_i) = \begin{cases} -1 & \text{if } \mathcal{L}_{a_i}[|s|_{a_i} + 1] > \mathcal{L}[|s| + 1] \\ 0 & \text{if } \mathcal{L}_{a_i}[|s|_{a_i} + 1] = \mathcal{L}[|s| + 1] \\ 1 & \text{if } \mathcal{L}_{a_i}[|s|_{a_i} + 1] < \mathcal{L}[|s| + 1] \end{cases}, \quad (1)$$

because we look for the minimal value over the singleton interval $\mathcal{L}[|s| + 1..|s| + 2]$, and for $c \neq a_i$ we have

$$b_c(s \cdot a_i) = \begin{cases} -1 & \text{if } \mathcal{L}_c[|s|_c] > \mathcal{L}[|s| + 1] \\ 0 & \text{if } \mathcal{L}_c[|s|_c] = \mathcal{L}[|s| + 1] \\ b_c(s) & \text{if } \mathcal{L}_c[|s|_c] < \mathcal{L}[|s| + 1] \end{cases}, \quad (2)$$

according to the relation of $\mathcal{L}[|s| + 1]$ to the minimal value in $\mathcal{L}[\ell_c(s) + 1..|s| + 2]$.

Now consider prefix consistency. The extension of s with a_i adds one new pair of consecutive occurrences of a_i 's, which satisfies the pair constraint if and only if $b_{a_i}(s) = 0$. The partial pair constraints of $s \cdot a_i$ can be checked using Lemma 8.

The computation of $b(s \cdot a_i)$ requires the verification of a separate condition for each $b_c(s \cdot a_i)$ for each $c \in \Sigma$, hence it could be done in time $O(\sigma)$. On the other hand, the computation of $p(s \cdot a_i)$ can be done in a constant time. \square

The structure of the automaton $\mathcal{A}_{\mathcal{L}}$ produced by the algorithm is as follows. Each state v of $\mathcal{A}_{\mathcal{L}}$ corresponds to a unique pair (p_v, b_v) and represents the set of strings $S_v = \{s : p(s) = p_v \wedge b(s) = b_v\}$. For a pair of states $v_1, v_2 \in \mathcal{A}_{\mathcal{L}}$ there exists a transition $v_1 \rightarrow v_2$ labelled with c if $S_{v_2} = S_{v_1} \cdot c$ and for each $s \in S_{v_1}$ sc is consistent with \mathcal{L} (where $A \cdot c$ denotes appending a character c to each element of the set A). In such a case (p_{v_2}, b_{v_2}) are given by the equations (1) and (2).

Note that if for a string s consistent with \mathcal{L} $b(s \cdot c)$ contains -1 , then the state v representing s can not have an outgoing transition labelled with c . Therefore, for any s consistent with \mathcal{L} , $b(s)$ can be represented as a bit vector (i.e. contain only binary values).

Observe that the empty string ε and all single characters $c \in \Sigma$ are consistent with \mathcal{L} . Hence, we can construct the set $W_{\mathcal{L}}$ and the automaton $\mathcal{A}_{\mathcal{L}}$ by iterative extension of strings consistent with \mathcal{L} . To construct $\mathcal{A}_{\mathcal{L}}$ we iterate through sets of states corresponding to strings of length $k = 1, \dots, n - 1$, i.e.

$$\mathcal{P}_k = \left\{ v \in \mathcal{A}_{\mathcal{L}} : \forall_{w \in S_v} |w| = k \right\},$$

and for each state $v \in \mathcal{P}_k$ we check the existence of a transition $v \rightarrow v_1$. All states corresponding to the sets of strings of length $k + 1$ consistent with \mathcal{L} form the set \mathcal{P}_{k+1} .

Observe that for any $w \in W_{\mathcal{L}}$ we have $p(w) = p(L)$ and $b(w) = (0, 0, \dots)$ (for each c $\mathcal{L}_c[|w|_c] = \mathcal{L}[|w|] = -2$). Therefore the final state v_f of $\mathcal{A}_{\mathcal{L}}$ is unique.

Now we are ready to discuss the time and space complexity of our solution. The number of states of $\mathcal{A}_{\mathcal{L}}$ is bounded by the number all possible pairs (p, b) of Parikh vectors and bit vectors. The number of all possible bit vectors is bounded by 2^σ and the number of all possible Parikh vectors reaches its maximum when the number of occurrences of all characters are equal. Moreover, we need $O(\sigma)$ space to store each state of $\mathcal{A}_{\mathcal{L}}$. Therefore, the space complexity of presented algorithm is $O(\sigma 2^\sigma (\frac{n}{\sigma} + 1)^\sigma)$.

To construct an automaton $\mathcal{A}_{\mathcal{L}}$ returned by the algorithm we need to check for each state $v \in \mathcal{A}_{\mathcal{L}}$ up to σ possible transition. Validation of a single transition requires $O(\sigma)$ time. This, together with the bound for the number of all states, gives us the time complexity $O(\sigma^2 2^\sigma (\frac{n}{\sigma} + 1)^\sigma)$.

The above discussion constitutes a proof of Theorem 5 in Section 8.

Remark 1. The above presented algorithm works correctly also for binary alphabet, however its time and space complexity is worse than the complexity of Algorithm 1.

The following examples illustrate the operation of reconstruction algorithm described above.

Example 4. Let us recall an integer array $\mathcal{L}[1..7] = [1, 4, 0, 2, 1, 1, 3]$ considered in Example 3. Using the procedure described above we will try to construct a finite deterministic automaton $\mathcal{A}_{\mathcal{L}}$ accepting the set of strings $W_{\mathcal{L}} = \{w \in \{a, b\}^7 : \text{LCP}_{\text{IBWT}(w)} = \mathcal{L}\}$.

First we transform \mathcal{L} into $\mathcal{L}[0..8] = [-1, 1, 4, 0, 2, 1, 3, -2]$ and compute character sequences $\mathcal{L}_a[0..4] = [-1, 0, 3, -2]$ and $\mathcal{L}_b[0..5] = [-1, 1, 0, 2, -2]$. The structure of $\mathcal{A}_{\mathcal{L}}$ is depicted on Figure 9.

We start with the automaton $\mathcal{A}_{\mathcal{L}}$ consisting of a single initial node $v_{(0)}$ represented by a pair $(p_0, b_0) = ([0, 0], [0, 0])$ and contained in the set \mathcal{P}_0 . Next, we iterate over all sets \mathcal{P}_k for $k = 0, \dots, n - 1$ and check for a possible extensions of each state $v \in \mathcal{P}_k$.

$\mathbf{v}_{(0)}$: By (1) and (2), $p(a) = [1, 0]$, $b(a) = [1, 0]$, $p(b) = [0, 1]$ and $b(b) = [0, 1]$. Since neither $b(a)$ nor $b(b)$ contain -1 , both strings are consistent with \mathcal{L} . Hence we create states $v_{(1)}$, $v_{(2)}$ and the transitions $v_{(0)} \rightarrow v_{(1)}$ and $v_{(0)} \rightarrow v_{(2)}$ labelled with a and b respectively.

$\mathbf{v}_{(1)}$: By (1) and (2) we have $p(aa) = [2, 0]$, $b(aa) = [1, 0]$, but $\mathcal{L}_a[2] = 3 \neq 1 = \mathcal{L}[2]$. Hence, due to pair constraint violation aa is not consistent with \mathcal{L} . On the other hand, $p(ab) = [1, 1]$, $b(ab) = [1, 1]$ and it is the first occurrence of b , hence we create a new state $v_{(3)}$ and a transition $v_{(1)} \rightarrow v_{(3)}$ labelled with b .

$\mathbf{v}_{(2)}$: We have $p(ba) = [1, 1]$, $b(ba) = [1, 0]$ and it is the first occurrence of a , hence we create a new state $v_{(4)}$ and a transition $v_{(2)} \rightarrow v_{(4)}$ labelled with a . We have $p(bb) = [0, 2]$, $b(bb) = [0, 1]$ and $\mathcal{L}_b[1] = 1 = \mathcal{L}[1]$, hence we create a new state $v_{(5)}$ and a transition $v_{(2)} \rightarrow v_{(5)}$ labelled with b .

$\mathbf{v}_{(3)}$: We have $p(aba) = [2, 1]$ and $b(aba) = [-1, -1]$. Moreover we have $p(abb) = [2, 1]$ and $b(abb) = [0, -1]$. Therefore, both aba and abb are not consistent with \mathcal{L} and $v_{(3)}$ has no valid extension. Due to that we remove states $v_{(3)}$ and $v_{(1)}$ (for which $v_{(3)}$ is the only successor) from $\mathcal{A}_{\mathcal{L}}$.

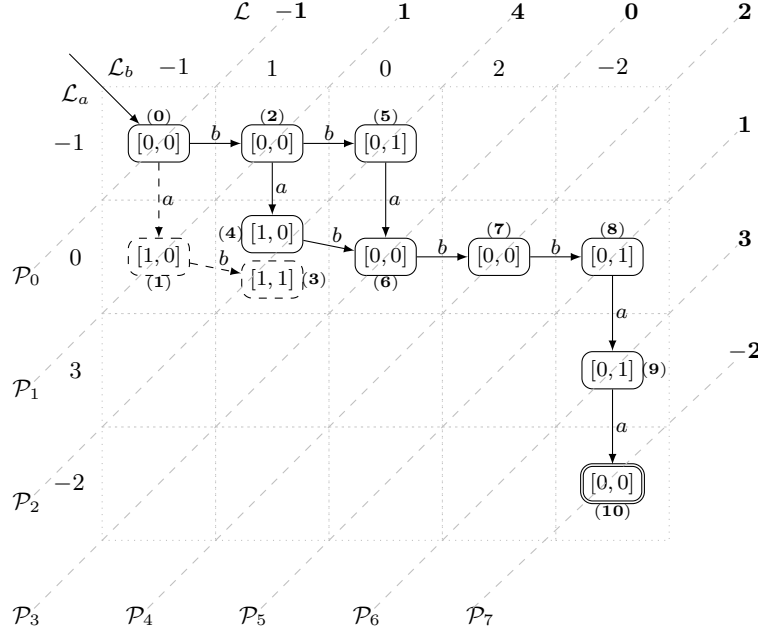


Fig. 9. The finite deterministic automaton $\mathcal{A}_{\mathcal{L}}$ constructed for LCP array $\mathcal{L} = [1, 4, 0, 2, 1, 3]$. Squares containing nodes represent Parikh vectors, i.e. the square in i -th row and j -th column represents the vector $p = (i - 1, j - 1)$. \mathcal{P}_k denotes the sets of nodes representing all strings of length k consistent with \mathcal{L} . The temporarily created states with no valid extension, which are not included in $\mathcal{A}_{\mathcal{L}}$, were marked with dashed lines. We have two possible paths leading from the initial to the final state corresponding to strings $babbbaa$ and $bbabbaa$.

$\mathbf{v}_{(4)}$: We have $p(baa) = [2, 1]$ and $b(baa) = [-1, -1]$, hence baa is not consistent with \mathcal{L} . On the other hand, we have $p(bab) = [1, 2]$, $b(bab) = [0, 0]$ and $\mathcal{L}_b[1] = 1 = \mathcal{L}[1]$, hence we create a new state $v_{(6)}$ and a transition $v_{(4)} \rightarrow v_{(6)}$ labelled with b .

$\mathbf{v}_{(5)}$: We have $p(bbb) = [1, 3]$ and $b(bbb) = [0, -1]$, hence bbb is not consistent with \mathcal{L} . On the other hand, we have $p(bba) = [1, 2]$, $b(bba) = [0, 0]$ and it is the first occurrence of a , hence we add a transition $v_{(5)} \rightarrow v_{(6)}$ labelled with a .

$\mathbf{v}_{(6)}$: Notice that $v_{(6)}$ represents the set of strings $S_{v_{(6)}} = \{bab, bba\}$. We have $p(S_{v_{(6)}} \cdot a) = [2, 2]$ and $b(S_{v_{(6)}} \cdot a) = [-1, 0]$, hence neither $baba$ nor $bbaa$ is consistent with \mathcal{L} . On the other hand, we have $p(S_{v_{(6)}} \cdot b) = [1, 3]$, $b(S_{v_{(6)}} \cdot b) = [0, 0]$ and $\mathcal{L}_b[2] = 0 = \mathcal{L}[3]$, hence we create a new state $v_{(7)}$ and a transition $v_{(6)} \rightarrow v_{(7)}$ labelled with b .

$\mathbf{v}_{(7)}$: Notice that $v_{(7)}$ represents the set of strings $S_{v_{(7)}} = \{babbb, bbab\}$. We have $p(S_{v_{(7)}} \cdot a) = [2, 3]$ and $b(S_{v_{(7)}} \cdot a) = [-1, -1]$, hence neither $babba$ nor $bbaba$ is consistent with \mathcal{L} . On the other hand, we have $p(S_{v_{(7)}} \cdot b) = [1, 4]$, $b(S_{v_{(7)}} \cdot b) = [0, 1]$ and $\mathcal{L}_b[3] = 2 = \mathcal{L}[4]$, hence we create a new state $v_{(8)}$ and a transition $v_{(7)} \rightarrow v_{(8)}$ labelled with b .

$\mathbf{v}_{(8)}$: Note that for each $s \in S_{v_{(8)}}$ we have $|s|_b = |\mathcal{L}|_b$, hence $s \cdot b$ is not consistent with \mathcal{L} . On the other hand, we have $p(S_{v_{(8)}} \cdot a) = [2, 4]$, $b(S_{v_{(8)}} \cdot a) = [0, 1]$ and $\mathcal{L}_a[1] = 0 = \mathcal{L}[3]$, hence we create a new state $v_{(9)}$ and a transition $v_{(8)} \rightarrow v_{(9)}$ labelled with a .

$\mathbf{v}_{(9)}$: Similarly as in the case of $v_{(8)}$, for each $s \in S_{v_{(9)}}$ we have $|s|_b = |\mathcal{L}|_b$, hence $s \cdot b$ is not consistent with \mathcal{L} . On the other hand, we have $p(S_{v_{(9)}} \cdot a) = [3, 4]$, $b(S_{v_{(9)}} \cdot b) = [0, 0]$ and $\mathcal{L}_a[2] = 3 = \mathcal{L}[6]$, hence we create a new state $v_{(10)}$ and a transition $v_{(9)} \rightarrow v_{(10)}$ labelled with a .

Finally, after computing $\mathcal{A}_{\mathcal{L}}$ and backtracking all the paths from $v_{(0)}$ to $v_{(10)}$ we obtain a set $W_{\mathcal{L}} = \{babbbbaa, bbabbaa\}$ (compare this to the result in Example 3).

Example 5. Let us consider an LCP array $\mathcal{L}[1..6] = [1, 0, 1, 0, 2]$. Using the CSSILA algorithm, we will try to reconstruct a set of strings $W_{\mathcal{L}} \subseteq \{a, b, c\}^6$, such that for each $w \in W_{\mathcal{L}}$ we have $\text{LCP}_{\text{IBWT}(w)} = \mathcal{L}$. Looking at the structure of \mathcal{L} we conclude that if there exist a solution, it must satisfy $|w|_a = |w|_b = |w|_c = 2$ for any $w \in W_{\mathcal{L}}$.

First we transform \mathcal{L} into $\mathcal{L}[0..7] = [-1, 1, 0, 1, 0, 2, -2]$ and compute character sequences $\mathcal{L}_a[0..3] = [-1, 0, -2]$, $\mathcal{L}_b[0..3] = [-1, 0, -2]$ and $\mathcal{L}_c[0..3] = [-1, 1, -2]$. The structure of the complete automaton $\mathcal{A}_{\mathcal{L}}$ is depicted on Figure 9.

We start with the automaton $\mathcal{A}_{\mathcal{L}}$ consisting of a single initial node $v_{(0)}$ represented by a pair $(p_0, b_0) = ([0, 0, 0], [0, 0, 0])$ and contained in the set \mathcal{P}_0 . Next, we iterate over all sets \mathcal{P}_k for $k = 0, \dots, n-1$ and check for a possible extensions of each state $v \in \mathcal{P}_k$. In all cases below we do not consider the obvious inconsistency of $s \cdot c$ for $|s|_c = |\mathcal{L}|_c$.

$\mathbf{v}_{(0)}$: By (1) and (2), $p(a) = [1, 0, 0]$, $b(a) = [1, 0, 0]$, $p(b) = [0, 1, 0]$, $b(b) = [0, 1, 0]$, $p(c) = [0, 0, 1]$, $b(c) = [0, 0, 0]$ and those are the first occurrences of each character. Since none of $b(a)$, $b(b)$ and $b(c)$ contain -1 , all singleton strings are consistent with \mathcal{L} . Hence we create states $v_{(1)}$, $v_{(2)}$ and $v_{(3)}$ and the transitions $v_{(0)} \rightarrow v_{(1)}$, $v_{(0)} \rightarrow v_{(2)}$ and $v_{(0)} \rightarrow v_{(3)}$ labelled with a , b and c respectively.

$\mathbf{v}_{(1)}$: We have $p(aa) = [2, 0, 0]$, $b(aa) = [1, 0, 0]$ but $\mathcal{L}_a[1] = 0 \neq 1 = \mathcal{L}[1]$, and $p(ac) = [1, 0, 1]$ and $b(ac) = [0, 0, -1]$, hence there is no consistent extension of $v_{(1)}$ with a and c . On the other hand, $p(ab) = [1, 1, 0]$, $b(ab) = [0, 0, 0]$ and it is the first occurrence of b . Hence we create state $v_{(4)}$ and the transition $v_{(1)} \rightarrow v_{(4)}$, labelled with b .

$\mathbf{v}_{(2)}$: We have $p(ba) = [1, 1, 0]$, $b(ba) = [0, 0, 0]$ and it is the first occurrence of a . Hence we create the transition $v_{(2)} \rightarrow v_{(4)}$, labelled with a . On the other hand, we have $p(bb) = [0, 2, 0]$, $b(bb) = [0, 1, 0]$, but $\mathcal{L}_b[1] = 0 \neq 1 = \mathcal{L}[1]$, and $p(bc) = [0, 2, 0]$, $b(bc) = [0, 0, -1]$, hence there is no consistent extension of $v_{(2)}$ with b and c .

$\mathbf{v}_{(3)}$: We have $p(ca) = [1, 0, 1]$, $b(bb) = [0, 0, -1]$, and $p(cb) = [0, 1, 1]$, $b(cb) = [0, 0, -1]$, hence there is no consistent extension of $v_{(3)}$ with a and b . On the other hand, we have $p(cc) = [0, 0, 2]$, $b(bb) = [0, 0, 1]$ and $\mathcal{L}_c[1] = 1 = \mathcal{L}[1]$, hence we create the state $v_{(5)}$ and the transition $v_{(3)} \rightarrow v_{(5)}$, labelled with c .

Summing up for \mathcal{P}_1 , ab , ba and cc are consistent with \mathcal{L} , while aa , ac , bb , bc , ca and cb are not.

$\mathbf{v}_{(4)}$: We have $p(S_{v_{(4)}} \cdot a) = [1, 2, 0]$, $b(S_{v_{(4)}} \cdot a) = [1, 0, 0]$ and $\mathcal{L}_a[1] = 0 = \mathcal{L}[2]$, hence we create the state $v_{(6)}$ and the transition $v_{(4)} \rightarrow v_{(6)}$, labelled with a . We have $p(S_{v_{(4)}} \cdot b) = [1, 2, 0]$, $b(S_{v_{(4)}} \cdot b) = [0, 1, 0]$ and $\mathcal{L}_b[1] = 0 = \mathcal{L}[2]$, hence we create the state $v_{(7)}$ and the transition $v_{(4)} \rightarrow v_{(7)}$, labelled with b . We have $p(S_{v_{(4)}} \cdot c) = [1, 1, 1]$, $b(S_{v_{(4)}} \cdot c) = [0, 0, 0]$ and it is the first occurrence of c , hence we create the state $v_{(8)}$ and the transition $v_{(4)} \rightarrow v_{(8)}$, labelled with c .

$\mathbf{v}_{(5)}$: We have $p(S_{v_{(5)}} \cdot a) = [1, 0, 2]$, $b(S_{v_{(5)}} \cdot a) = [1, 0, 1]$ and it is the first occurrence of a , hence we create the state $v_{(9)}$ and the transition $v_{(5)} \rightarrow v_{(9)}$, labelled with a . We have $p(S_{v_{(5)}} \cdot b) = [0, 1, 2]$, $b(S_{v_{(5)}} \cdot b) = [0, 1, 1]$ and it is the first occurrence of b , hence we create the state $v_{(10)}$ and the transition $v_{(5)} \rightarrow v_{(10)}$, labelled with b .

Summing up for \mathcal{P}_2 , aba , abb , abc , baa , bab , bac , cca and ccb are consistent with \mathcal{L} .

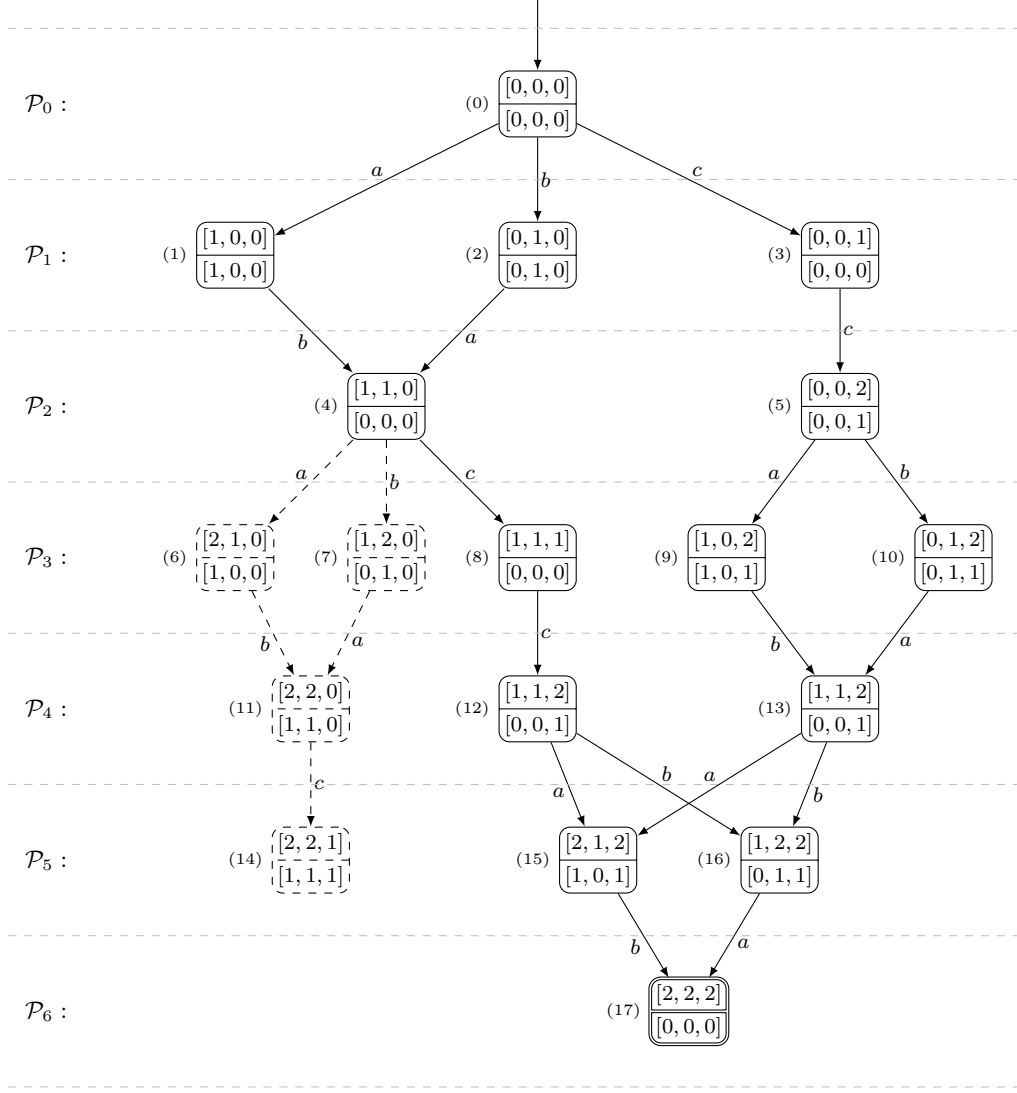


Fig. 10. The finite deterministic automaton $\mathcal{A}_{\mathcal{L}}$ constructed for LCP array $\mathcal{L} = [1, 0, 1, 0, 2]$. The upper part of each state v contains the Parikh vector p_v , while the lower part the bit vector b_v . \mathcal{P}_k denotes the sets of nodes representing all strings of length k consistent with \mathcal{L} . The temporarily created states with no valid extension, which are not included in $\mathcal{A}_{\mathcal{L}}$, were marked with dashed lines. We have eight possible paths leading from the initial to the final state corresponding to strings $abccab$, $abccba$, $baccab$, $baccba$, $ccabab$, $ccabba$, $ccbaab$ and $ccbaba$.

$\mathbf{v}_{(6)}$: We have $p(S_{v_{(6)}} \cdot b) = [2, 2, 0]$, $b(S_{v_{(6)}} \cdot b) = [1, 1, 0]$ and $\mathcal{L}_b[1] = 0 = \mathcal{L}[2]$, hence we create the state $v_{(11)}$ and the transition $v_{(6)} \rightarrow v_{(11)}$, labelled with b . On the other hand, we have $p(S_{v_{(6)}} \cdot c) = [2, 1, 1]$ and $b(S_{v_{(6)}} \cdot c) = [1, 0, -1]$, hence there is no consistent extension of $v_{(6)}$ with c .

$\mathbf{v}_{(7)}$: We have $p(S_{v_{(7)}} \cdot a) = [2, 2, 0]$, $b(S_{v_{(7)}} \cdot a) = [1, 1, 0]$ and $\mathcal{L}_a[1] = 0 = \mathcal{L}[2]$, hence we create the transition $v_{(7)} \rightarrow v_{(11)}$, labelled with a . On the other hand, we have $p(S_{v_{(6)}} \cdot c) = [1, 2, 1]$ and $b(S_{v_{(7)}} \cdot c) = [0, 0, -1]$, hence there is no consistent extension of $v_{(7)}$ with c .

$\mathbf{v}_{(8)}$: We have $p(S_{v_{(8)}} \cdot c) = [1, 1, 2]$, $b(S_{v_{(6)}} \cdot c) = [0, 0, 1]$ and $\mathcal{L}_c[1] = 1 = \mathcal{L}[3]$, hence we create the state $v_{(12)}$ and the transition $v_{(8)} \rightarrow v_{(12)}$, labelled with c . On the other hand we have $p(S_{v_{(8)}} \cdot a) = [2, 1, 1]$ and $b(S_{v_{(6)}} \cdot a) = [0, 0, -1]$, and $p(S_{v_{(8)}} \cdot b) = [1, 2, 1]$, $b(S_{v_{(6)}} \cdot b) = [0, 0, -1]$, hence there is no consistent extension of $v_{(8)}$ with a and b .

$\mathbf{v}_{(9)}$: We have $p(S_{v_{(9)}} \cdot b) = [1, 1, 2]$, $b(S_{v_{(9)}} \cdot b) = [0, 0, 1]$ and it is the first occurrence of b , hence we create the state $v_{(13)}$ and the transition $v_{(9)} \rightarrow v_{(13)}$, labelled with b . On the other hand, we have $p(S_{v_{(9)}} \cdot a) = [2, 0, 2]$, $b(S_{v_{(9)}} \cdot a) = [1, 0, 1]$, but $\mathcal{L}_a[1] = 0 \neq 1 = \mathcal{L}[3]$, hence there is no consistent extension of $v_{(9)}$ with a .

$\mathbf{v}_{(10)}$: We have $p(S_{v_{(10)}} \cdot a) = [1, 1, 2]$, $b(S_{v_{(10)}} \cdot a) = [0, 0, 1]$ and it is the first occurrence of a , hence we create the transition $v_{(10)} \rightarrow v_{(13)}$, labelled with a . On the other hand, we have $p(S_{v_{(10)}} \cdot b) = [0, 2, 2]$, $b(S_{v_{(10)}} \cdot b) = [0, 1, 1]$, but $\mathcal{L}_b[1] = 0 \neq 1 = \mathcal{L}[3]$, hence there is no consistent extension of $v_{(10)}$ with b .

Summing up for \mathcal{P}_3 , *abab*, *abba*, *abcc*, *baab*, *baba*, *bacc*, *ccab* and *ccba* are consistent with \mathcal{L} , while *abac*, *abbc*, *abca*, *abcb*, *baac*, *babc*, *bacab*, *ccaa* and *ccbb* are not.

$\mathbf{v}_{(11)}$: We have $p(S_{v_{(11)}} \cdot c) = [2, 2, 1]$, $b(S_{v_{(11)}} \cdot c) = [1, 1, 1]$ and it is the first occurrence of c , hence we create the state $v_{(14)}$ and the transition $v_{(11)} \rightarrow v_{(12)}$, labelled with c .

$\mathbf{v}_{(12)}$: We have $p(S_{v_{(12)}} \cdot a) = [2, 1, 2]$, $b(S_{v_{(12)}} \cdot a) = [1, 0, 1]$ and $\mathcal{L}_a[1] = 1 = \mathcal{L}[4]$, hence we create the state $v_{(15)}$ and the transition $v_{(12)} \rightarrow v_{(15)}$, labelled with a . Similarly, we have $p(S_{v_{(12)}} \cdot b) = [1, 2, 2]$, $b(S_{v_{(12)}} \cdot b) = [0, 1, 1]$ and $\mathcal{L}_b[1] = 1 = \mathcal{L}[4]$, hence we create the state $v_{(16)}$ and the transition $v_{(12)} \rightarrow v_{(16)}$, labelled with b .

$\mathbf{v}_{(13)}$: We have $p(S_{v_{(13)}} \cdot a) = [2, 1, 2]$, $b(S_{v_{(13)}} \cdot a) = [1, 0, 1]$ and $\mathcal{L}_a[1] = 1 = \mathcal{L}[4]$, hence we the transition $v_{(13)} \rightarrow v_{(15)}$, labelled with a . Similarly, we have $p(S_{v_{(13)}} \cdot b) = [1, 2, 2]$, $b(S_{v_{(13)}} \cdot b) = [0, 1, 1]$ and $\mathcal{L}_b[1] = 1 = \mathcal{L}[4]$, hence we create the transition $v_{(13)} \rightarrow v_{(16)}$, labelled with b .

Summing up for \mathcal{P}_4 , *ababc*, *abbac*, *abcca*, *abccb*, *baabc*, *babac*, *bacca*, *baccb*, *ccaba*, *ccabb*, *ccbaa* and *ccbab* are consistent with \mathcal{L} .

$\mathbf{v}_{(14)}$: We have $p(S_{v_{(14)}} \cdot c) = [2, 2, 2]$, $b(S_{v_{(14)}} \cdot c) = [0, 0, 0]$, but $\mathcal{L}_c[1] = 1 \neq 2 = \mathcal{L}[5]$, hence there is no consistent extension of $v_{(14)}$ with c .

Since $v_{(14)}$ has no extension and is not a part of the solution, it should be removed from $\mathcal{A}_{\mathcal{L}}$. This implies also removing $v_{(11)}$, which has $v_{(14)}$ as its only extension, and further removing $v_{(6)}$ and $v_{(7)}$ both having $v_{(11)}$ as their only extension.

$\mathbf{v}_{(15)}$: We have $p(S_{v_{(15)}} \cdot b) = [2, 2, 2]$, $b(S_{v_{(15)}} \cdot b) = [0, 0, 0]$ and $\mathcal{L}_b[1] = 1 = \mathcal{L}[4]$, hence we create the state $v_{(17)}$ and the transition $v_{(15)} \rightarrow v_{(17)}$, labelled with b .

$\mathbf{v}_{(16)}$: We have $p(S_{v_{(16)}} \cdot a) = [2, 2, 2]$, $b(S_{v_{(16)}} \cdot a) = [0, 0, 0]$ and $\mathcal{L}_a[1] = 1 = \mathcal{L}[4]$, hence we create the transition $v_{(16)} \rightarrow v_{(17)}$, labelled with a .

Summing up, *abccab*, *abccba*, *baccab*, *baccba*, *ccabab*, *ccabba*, *ccbaab* and *ccbaba* are consistent with \mathcal{L} , while *ababcc*, *abbacc*, *baabcc* and *babacc* are not.

Finally, after computing $\mathcal{A}_{\mathcal{L}}$ we can recover the set

$$W_{\mathcal{L}} = \{abccab, abccba, baccab, baccba, ccabab, ccabba, ccbaab, ccbaba\}$$

by backtracking all paths leading from the initial node $v_{(0)}$ to the final node $v_{(17)}$.

Note that to list the strings which are inconsistent with \mathcal{L} for each hyperplane \mathcal{P}_k we consider only those having prefixes of length $k-1$, which are consistent with \mathcal{L} . If we skip this requirement, we can produce more examples of strings inconsistent with \mathcal{L} .

E BCSILA to CCEC: An Example

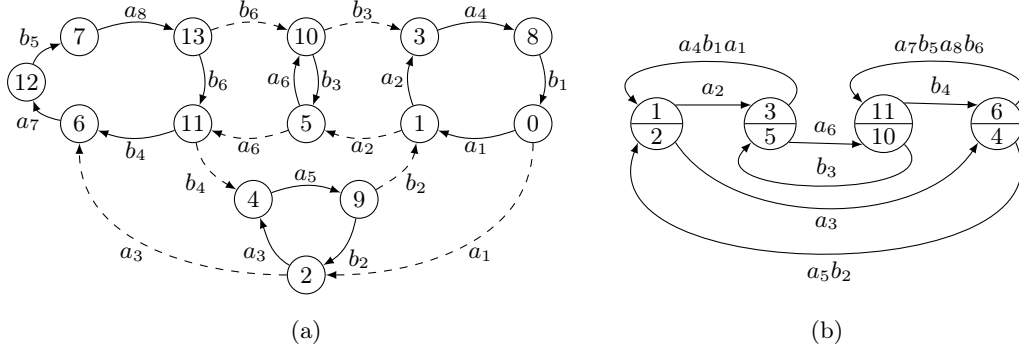


Fig. 11. The graphs G_V (a) and \tilde{G}_V (b) for $V = b[ab][aabb]baa[ab]aa$, which is the BWT with swaps produced from the LCP array $\mathcal{L} = [2, 5, 1, 4, 3, 4, 2, 0, 3, 2, 5, 3, 1]$. The solid edges in G_V are the edges of G_v for $v = babaabbbbaaaba$. The graph \tilde{G}_V is the CCEC instance derived from the BCSILA instance \mathcal{L} . In the initial state, all vertices are in the straight state, so that the cycles in \tilde{G}_V correspond to the cycles in G_v . The only non-singleton partition in \tilde{G}_v is $\{3/5, 6/4\}$ corresponding to the only swap interval of length more than two in V .

F Identification of Swap Cores

As descibed in Section 6, the BCSILA instance derived from a CCEC instance has the strings $ba^k b$, $k \in [1..m]$, as desired swap cores. We will next systematically inspect all other strings to identify all other (undesirable) swap cores. Recall that an interval $[i..j]$ in V is a swap interval if and only if the following conditions hold:

1. $[i..j]$ is an x -interval for a string x such that either $occ(axa) = occ(bxb) = occ(x)/2$ or $occ(axb) = occ(bxa) = occ(x)/2$, where $occ(y)$ is the number of occurrences of y in W , and
2. $LCP_W[i + 1..k] = LCP_W[k + 1..j]$, where $k = (i + j)/2$.

Notice that if $occ(x) = j - i = 2$, the second condition is trivially true.

Let us start with unary strings. First, b , bb and a^k for $k < m + 2n$ are not swap cores because they are preceded and succeeded by a more often than by b . We can also eliminate all other unary strings since they occur at most once. We also note that any string beginning (ending) with bb cannot be a swap core because it is always preceded (succeeded) by a . Let us then consider strings x of the following forms:

- $x = ba^k$. If $k < m + 2n - 1$, $occ(xa) > occ(xb)$, and if $k \geq m + 2n$, $occ(x) \leq 1$. In either case, x is not a swap core. On the other hand, $x = ba^{m+2n-1}$ is always a swap core with two occurrences.
- $x = a^k b$. This case is symmetric to the one above except we cannot be certain whether $x = a^{m+2n-1} b$ is a swap core or not since the characters following the two occurrences of x are not fully determined. However, we count x as a potential swap core.
- $x = a^k ba^k$ and $x = a^k bba^k$. If $k > m$, we have $occ(x) = 0$, and if $k < m$, we have $occ(ax) > occ(bx)$. If $k = m$, then x is a swap core and $occ(x) = 2$.
- $x = a^k ba^h$ and $x = a^k bba^h$ for $k < h$. If $k > m$, we have $occ(x) = 0$ and if $h \leq m$, we have $occ(ax) > occ(bx)$. If $k \leq m < h$, x is obviously not a swap core if $occ(x) < 2$ but also not if $occ(x) > 2$ because then we must have $occ(xa) > occ(xb)$. On the other hand, if $k \leq m < h$ and $occ(x) = 2$, then x might be a swap core.

- $x = a^kba^h$ and $x = a^kbba^h$ for $k > h$. This is symmetric to the case above.
- $x = ba^kba^h$, $x = ba^kbba^h$, $x = a^hba^kb$ and $x = a^hbba^kb$. If $k > m$, $occ(x) \leq 1$. If $k \leq m$, every occurrence of x is either preceded (the first two cases) or succeeded (the latter two cases) by the same character. Thus x is never a swap core.
- $x = a^kba^iba^h$ and $x = a^kbba^ibba^h$ for $i \in [1..m]$. Obviously, x is not a swap core if $occ(x) < 2$ but also not if $occ(x) > 2$ because then $occ(xa) > occ(xb)$. If $occ(x) = 2$ then x may or may not be a swap core.

Any string not mentioned above either does not occur at all or contains a substring of the form ba^kb for $k > m$ and occurs once.

Notice that each of the potential extra swap cores has exactly two occurrences.